

AIX Version 7.2

*Technical Reference: Kernel and  
Subsystems, Volume 1*

**IBM**



AIX Version 7.2

*Technical Reference: Kernel and  
Subsystems, Volume 1*

**IBM**

**Note**

Before using this information and the product it supports, read the information in "Notices" on page 695.

This edition applies to AIX Version 7.2 and to all subsequent releases and modifications until otherwise indicated in new editions.

© **Copyright IBM Corporation 2015, 2017.**

US Government Users Restricted Rights – Use, duplication or disclosure restricted by GSA ADP Schedule Contract with IBM Corp.

---

# Contents

<b>About this document</b> . . . . .	<b>v</b>
Highlighting . . . . .	v
Case sensitivity in AIX . . . . .	v
ISO 9000. . . . .	v

## **Technical Reference: Kernel and Subsystems, Volume 1** . . . . . **1**

What's new in Technical Reference: Kernel and Subsystems, Volume 1 . . . . .	1
Kernel Services . . . . .	1
a . . . . .	1
b . . . . .	25
c . . . . .	37
d . . . . .	61
e . . . . .	112
f . . . . .	140
g . . . . .	181
h . . . . .	201
i . . . . .	207
k . . . . .	240
l . . . . .	328
m . . . . .	356
n . . . . .	377
p . . . . .	399
q . . . . .	428
r . . . . .	429
s . . . . .	462
t . . . . .	481
u . . . . .	504
v . . . . .	532
w . . . . .	582
x . . . . .	591
Device Driver Operations . . . . .	614
Standard Parameters to Device Driver Entry Points . . . . .	614
buf Structure . . . . .	615
bufx Structure . . . . .	617
Character Lists Structure . . . . .	619
ddclose Device Driver Entry Point . . . . .	620
ddconfig Device Driver Entry Point . . . . .	621
dddump Device Driver Entry Point . . . . .	624
ddioctl Device Driver Entry Point . . . . .	626
ddmpx Device Driver Entry Point . . . . .	627
ddopen Device Driver Entry Point . . . . .	629
ddread Device Driver Entry Point . . . . .	631
ddrevoke Device Driver Entry Point . . . . .	632
ddselect Device Driver Entry Point . . . . .	633
ddstrategy Device Driver Entry Point . . . . .	635
ddwrite Device Driver Entry Point . . . . .	637
Select/Poll Logic for ddwrite and ddread Routines . . . . .	639
uio Structure . . . . .	639
Virtual File System Operations . . . . .	641
vfs_aclxctl Entry Point . . . . .	642
vfs_cntl Entry Point . . . . .	644

vfs_hold or vfs_unhold Kernel Service . . . . .	645
vfs_init Entry Point . . . . .	645
vfs_mount Entry Point . . . . .	646
vfs_root Entry Point . . . . .	648
vfs_search Kernel Service . . . . .	648
vfs_statfs Entry Point . . . . .	649
vfs_sync Entry Point . . . . .	650
vfs_umount Entry Point . . . . .	651
vfs_vget Entry Point . . . . .	652
vnop_access Entry Point . . . . .	653
vnop_close Entry Point . . . . .	655
vnop_create Entry Point . . . . .	656
vnop_create_attr Entry Point . . . . .	657
vnop_fclear Entry Point . . . . .	658
vnop_fid Entry Point . . . . .	659
vnop_finfo Entry Point . . . . .	659
vnop_fsync, vnop_fsync_range Entry Points . . . . .	660
vnop_ftrunc Entry Point . . . . .	662
vnop_getacl Entry Point . . . . .	662
vnop_getattr Entry Point . . . . .	663
vnop_getxacl Entry Point . . . . .	664
vnop_hold Entry Point . . . . .	665
vnop_ioctl Entry Point . . . . .	666
vnop_link Entry Point . . . . .	667
vnop_lockctl Entry Point . . . . .	668
vnop_lookup Entry Point . . . . .	670
vnop_map Entry Point . . . . .	671
vnop_map_lloff Entry Point . . . . .	672
vnop_memcntl Entry Point . . . . .	673
vnop_mkdir Entry Point . . . . .	674
vnop_mknod Entry Point . . . . .	675
vnop_open Entry Point . . . . .	675
vnop_rdw, vnop_rdw_attr Entry Points . . . . .	676
vnop_readdir Entry Point . . . . .	678
vnop_readdir_eofp Entry Point . . . . .	679
vnop_readlink Entry Point . . . . .	680
vnop_rele Entry Point . . . . .	681
vnop_remove Entry Point . . . . .	681
vnop_rename Entry Point . . . . .	682
vnop_revoke Entry Point . . . . .	683
vnop_rmdir Entry Point . . . . .	684
vnop_seek Entry Point . . . . .	685
vnop_select Entry Point . . . . .	686
vnop_setacl Entry Point . . . . .	687
vnop_setattr Entry Point . . . . .	688
vnop_setxacl Entry Point . . . . .	689
vnop_strategy Entry Point . . . . .	691
vnop_symlink Entry Point . . . . .	692
vnop_unmap Entry Point . . . . .	692

## **Notices** . . . . . **695**

Privacy policy considerations . . . . .	697
Trademarks . . . . .	697

## **Index** . . . . . **699**



---

## About this document

This topic collection is part of the six-volume technical reference set that provides information about system calls, kernel extension calls, and subroutines.

---

## Highlighting

The following highlighting conventions are used in this document:

<b>Bold</b>	Identifies commands, subroutines, keywords, files, structures, directories, and other items whose names are predefined by the system. Bold highlighting also identifies graphical objects, such as buttons, labels, and icons that the you select.
<i>Italics</i>	Identifies parameters for actual names or values that you supply.
Monospace	Identifies examples of specific data values, examples of text similar to what you might see displayed, examples of portions of program code similar to what you might write as a programmer, messages from the system, or text that you must type.

---

## Case sensitivity in AIX

Everything in the AIX® operating system is case sensitive, which means that it distinguishes between uppercase and lowercase letters. For example, you can use the **ls** command to list files. If you type **LS**, the system responds that the command is not found. Likewise, **FILEA**, **FiLea**, and **filea** are three distinct file names, even if they reside in the same directory. To avoid causing undesirable actions to be performed, always ensure that you use the correct case.

---

## ISO 9000

ISO 9000 registered quality systems were used in the development and manufacturing of this product.





---

## Technical Reference: Kernel and Subsystems, Volume 1

This topic collection provides information about kernel services, description of standard device driver entry points parameters, and the list of virtual file system operations.

The AIX operating system is designed to support The Open Group's Single UNIX Specification Version 3 (UNIX 03) for portability of operating systems based on the UNIX operating system. Many new interfaces, and some current ones, have been added or enhanced to meet this specification. To determine the correct way to develop a UNIX 03 portable application, see The Open Group's UNIX 03 specification on The UNIX System website (<http://www.unix.org>).

---

## What's new in Technical Reference: Kernel and Subsystems, Volume 1

Read about new or significantly changed information for the Technical Reference: Kernel and Subsystems, Volume 1 topic collection.

### How to see what's new or changed

In PDF files, you might see revision bars (|) in the left margin of new and changed information.

### January 2017

Added the `sleepx` and `in_localaddr` kernel services.

### December 2017

Added information about the `V_READMAKE` flag in the `vm_readp` kernel services.

### May 2016

Added information about the `vsx_enable` and the `vsx_disable` kernel services.

### December 2015

The following information is a summary of the updates made to this topic collection:

- Added information about the following kernel services:
  - “`fp_get_path` Kernel Service” on page 150
  - “`fskv_reg` Kernel Service” on page 175
  - “`fskv_unreg` Kernel Service” on page 179
  - “`nameToXfid()` Kernel Service” on page 377
  - “`TE_verify_reg` Kernel Service” on page 481
  - “`TE_verify_unreg` Kernel Service” on page 483
  - “`xfidToName()` Kernel Service” on page 591

---

## Kernel Services

The following kernel services begin with the with the letter a - x.

### a

The following kernel services begin with the with the letter a.

## **acct\_add\_LL or acct\_zero\_LL Kernel Service Purpose**

Increments counters for advanced accounting.

### **Syntax**

```
unsigned long long acct_add_LL(ptr, incr)
unsigned long long *ptr;
unsigned int incr;
```

```
unsigned long long acct_zero_LL(ptr)
unsigned long long *ptr;
```

### **Parameters**

<b>Item</b>	<b>Description</b>
<i>ptr</i>	Address of statistic to be incremented.
<i>incr</i>	Increment to be applied.

### **Description**

These kernel services are special atomic increment and clear services that are designed to allow machine-independent updating of **unsigned long long** values. The increment service only performs an increment if advanced accounting is enabled.

The **acct\_add\_LL** kernel service adds the value associated with the *incr* parameter to the 64-bit counter at the address designated by the *ptr* parameter. The **acct\_zero\_LL** kernel service atomically zeroes the 64-bit counter.

Both routines return the previous value of the 64-bit counter. This way, the **acct\_zero\_LL** kernel service can be used to atomically get the most recent value and set the counter to NULL. Because only delta statistics are reported each interval, this capability is required by interval accounting when the accounting record is being built for a report.

### **Execution Environment**

These kernel services can be called from either the interrupt environment or the process environment.

### **Return Values**

These subroutines return the previous value of the location designated by the *ptr* parameter.

#### **Related reference:**

“acct\_interval\_register or acct\_interval\_unregister Kernel Service” on page 4

“acct\_put Kernel Service” on page 6

## **acct\_get\_projid Kernel Service Purpose**

Gets the project identifier for the current process.

### **Syntax**

```
projid_t acct_get_projid(void)
```

## Description

The `acct_get_projid` kernel service returns the project identifier for the current process.

## Execution Environment

The `acct_get_projid` kernel service can be called from the process environment only.

## Return Values

The `acct_get_projid` kernel service returns the current project identifier.

### Related reference:

“`acct_put` Kernel Service” on page 6

## `acct_get_usage` Kernel Service Purpose

Allows kernel extensions to measure the resource utilization of transactions.

## Syntax

```
#include <sys/types.h>
#include <sys/aacct.h>

unsigned long long acct_get_usage(usage)
struct tusage *usage;
```

## Parameters

Item	Description
<i>usage</i>	Resource utilization structure.

## Description

This routine is used to measure the resource utilization of a client transaction, so that the cost of the transaction can be included within the accounting record that identifies the client transaction. This accounting record is then used for chargeback purposes.

The `acct_get_usage` kernel service is designed to be called twice: once at the start of a transaction and a second time at the end of a transaction. Each time that the routine is called, it returns the resource utilization for the calling thread from creation using the *usage* parameter. Therefore, this routine can be called multiple times to determine the resource utilization of a code fragment by subtracting start and end values.

The following macros are provided for manipulating the usage parameter:

### `TUSAGE_ZERO(TU)`

Initializes the `tusage` structure

### `TUSAGE_ADD(TU1, TU2)`

Adds `tusage` structures ( $T1 = T1 + T2$ )

### `TUSAGE_SUB(TU1, TU2)`

Subtracts `tusage` structures ( $T1 = T1 - T2$ )

The *usage* parameter provides thread-specific information, so the caller must ensure that this routine is called from the same thread context when measuring the utilization of a transaction. The return value identifies the calling thread context.

The `acct_get_usage` kernel service returns a token that identifies the calling context. This token can be logically compared with other tokens returned by this routine to ensure that start and stop invocations were made from the same thread. The scope of the token depends on the context of the calling program. If this routine is called under a pthread, then it returns a token representing the currently executing pthread. Otherwise, the `acct_get_usage` kernel service returns a token representing the currently executing kernel thread. In the former case, the token has process-wide scope; in the latter case, the token has system-wide scope.

## Execution Environment

The `acct_get_usage` kernel service can only be called from the process environment.

## Return Values

Upon successful completion, the `acct_get_usage` kernel service returns a token that identifies the calling thread context.

### Related reference:

“`acct_get_projid` Kernel Service” on page 2

“`acct_interval_register` or `acct_interval_unregister` Kernel Service”

“`acct_put` Kernel Service” on page 6

## `acct_interval_register` or `acct_interval_unregister` Kernel Service Purpose

Registers or unregisters an advanced accounting handler.

## Syntax

```
#include <sys/aacct.h>
```

```
int acct_interval_register(trid, cmds, handler, arg, reg_token, reg_name)
int trid;
int cmds;
int (*handler)(int trid, int cmds, void *arg);
void *arg;
unsigned long *reg_token;
char *reg_name;
```

```
int acct_interval_unregister(reg_token)
unsigned long reg_token;
```

## Parameters

Item	Description
<i>trid</i>	Transaction identifier
<i>cmds</i>	Invocations supported by the advanced accounting handler
<i>handler</i>	Function descriptor for the handler
<i>arg</i>	Identifies the instance of the kernel extension
<i>reg_token</i>	Token that is returned to caller naming the instance of the registration
<i>reg_name</i>	Identifies the transaction using a string

## Description

The `acct_interval_register` kernel service registers accounting records that are produced by the kernel extension with the advanced accounting subsystem. These accounting records are named through accounting transaction identifiers, which are provided by the caller. Transaction identifiers are persistent in nature, because they are used by report and analysis utilities to interpret transaction-specific accounting data. The transaction identifier is implicitly mapped to a template.

Transaction identifiers (and associated templates) used by AIX are defined in the **sys/aacct.h** file. Identifiers in the range of 0 – 127 are reserved for AIX. Vendors can choose any value in the range 128 – 256 for their accounting records. If two vendors choose the same value, report and analysis programs must reference other fields in the accounting record header to uniquely identify the source of the transaction; that way, they can apply the appropriate template. The *subproject* field (which specifies the command name of the logger) and *length* field can be used to identify the source of the transaction. Collisions are very unlikely to occur. The transaction identifier and the transaction name, which is provided by the *reg\_name* field, are presented to the system administrator. Vendors should choose representative names for their transactions. The maximum length of a transaction name is 15 bytes.

Administrators can enable and disable transactions, and thereby drive callouts to the kernel extension. A function descriptor for the advanced accounting handler is provided through the *handler* parameter. The interface of this handler is:

```
int handler(int trid, int cmd, void *arg);
```

The *trid* parameter is the transaction being acted on. The *cmd* parameter describes the action. The *arg* parameter is a value that was specified at registration for this particular instance of the handler. The *arg* parameter is specific to the kernel extension.

The following *cmd* values are supported:

Item	Description
ACCT_CMD_ENABLE	The transaction is being enabled; start collecting.
ACCT_CMD_DISABLE	The transaction is being disabled; stop collecting.
ACCT_CMD_INTERVAL	The system interval has expired; provide accounting data.
ACCT_CMD_FSWITCH	The active accounting file has changed; provide meta data.

The handler is invoked in the process environment from a dedicated kernel-only thread that is part of the advanced accounting subsystem. The kernel extension registers for the callouts that should be made by logically ORing *cmd* values. The *cmds* parameter to the **acct\_interval\_register** kernel service is provided for this purpose.

When a transaction is enabled, the kernel extension should allocate accounting structures and start collecting statistics. When a transaction is disabled, the kernel extension should quit collecting statistics and free accounting structures. If a transaction is not enabled, the kernel subsystem should not collect statistics for the transaction. The kernel extension relies on the callout mechanism to provide notification when a transaction is enabled. This way, accounting records that are not required for the report are not collected and the accounting overhead is minimized.

If the kernel extension registers for interval accounting, the extension is called when the system interval expires. The handler should record its data using the **acct\_put** kernel service and should reset its counters so that only delta statistics are produced in the next interval. The **acct\_zero\_LL** and **acct\_add\_LL** kernel services are provided so that statistics can be reported and zeroed atomically. When the system interval is disabled, the system automatically generates an interval callout to collect the last round of statistics.

The file switch callout is provided, so that subsystems can record accounting data in each accounting file. Most subsystems are not expected to use this option.

## Execution Environment

The **acct\_interval\_register** kernel service can be called from the process environment only.

The **acct\_interval\_unregister** kernel service can be called from either the interrupt environment or the process environment.

## Return Values

Upon successful completion, 0 is returned. If unsuccessful, **errno** is set to a value that explains the error.

### Related reference:

“acct\_add\_LL or acct\_zero\_LL Kernel Service” on page 2

“acct\_put Kernel Service”

## acct\_put Kernel Service

### Purpose

Writes an accounting record.

### Syntax

```
#include <sys/aacct.h>
```

```
void acct_put(trid, flags, projid, usage, trdata, tr_len);
int trid;
int flags;
projid_t projid;
struct tusage *usage;
void *trdata;
int tr_len;
```

### Parameters

Item	Description
<i>trid</i>	Transaction identifier.
<i>flags</i>	Flags associated with the transaction or the production of the transaction. The following value is defined:  <b>ACCT_PUT_DIRECT</b> Overrides aggregate transaction
<i>projid</i>	Project identifier, associated with the transaction, that identifies the billable entity. The following values are defined:  <b>PROJID_SYSTEM</b> This identifier is typically associated with system overhead and is often used for shared devices, such as disks and network adapters.  <b>PROJID_UNKNOWN</b> This identifier is used when the billable entity is unknown to the caller. In this case, the system calculates the project identifier using the project assignment policy specified by the system administrator.  <i>project identifier</i> If the project identifier is known, it should be specified.
<i>usage</i>	Identifies the resource usage values associated with the transaction.
<i>trdata</i>	Transaction-specific information.
<i>tr_len</i>	Size of the transaction-specific data in bytes.

### Description

The **acct\_put** kernel service provides accounting data to the advanced accounting subsystem. This service builds the accounting record header from its parameters and values associated with the calling context. The transaction-specific data specified by the caller is copied after the header. This data is internally buffered so that it can be written efficiently to the accounting data file some time later.

The *trid* parameter identifies the type of transaction that is being provided and implicitly identifies the format of the transaction-specific data. This identifier is included within the accounting header and is used by report and analysis commands to infer the right template that can interpret transaction-specific

data. Vendors are encouraged to document their transaction identifiers and record templates so that report and analysis tools can be produced to interpret this data.

Accounting transaction identifiers are defined in the following range:

Item	Description
0-127	AIX accounting transaction identifiers
128-255	Vendor accounting transaction identifiers

The **ACCT\_PUT\_DIRECT** flag is provided as an override to the aggregation of accounting records, which is an optional feature of the advanced accounting subsystem. By default, the system does not aggregate accounting data. Aggregation is designed to reduce the volume of data that is written to the accounting file. It is transparent to applications and middleware. When aggregation is enabled, the system throws out the transaction-specific data and produces statistics about the occurrence of the transaction and the aggregate resource utilization. The data is produced along project boundaries, so the ability to perform chargeback is not lost, although the data that is produced is different. Statistical information about the transaction is captured in the accounting file in lieu of the transaction.

Because aggregation might not be desirable in some cases, the **ACCT\_PUT\_DIRECT** flag is provided to override this feature. For example, because the significance of a transaction that describes the shared use of a disk is bound up in the transaction-specific data, the transaction cannot be effectively aggregated. The significance of the transaction is thrown out in the course of aggregation. In effect, the statistic has already been aggregated by the producer, so it should be written directly to the file instead of being aggregated again by the accounting subsystem.

The usage values pointed to by the *usage* parameter is calculated using the **acct\_get\_usage** kernel service. The *usage* parameter is optional. A value of NULL can be specified to signify no usage information. Aggregation uses this field to accumulate resource utilization. If this information is calculated for the transaction, it should be passed as a parameter to this routine, instead of just including it within the transaction-specific data section. The advanced accounting subsystem does not know the format of this section and cannot aggregate it. In such a case, this section would be thrown out when aggregation is enabled.

The *trdata* parameter contains the address of a buffer containing transaction-specific data, and the *tr\_len* parameter identifies the number of bytes in this buffer that should be copied to the accounting file. A maximum of 16 KB of data can be written.

## Execution Environment

The **acct\_put** kernel service can be started from either the process or interrupt environment. However, aggregation of the transaction is only supported when the **acct\_put** service is started from the process environment.

## Return Values

The **acct\_put** kernel service does not return a value.

### Related reference:

“**acct\_add\_LL** or **acct\_zero\_LL** Kernel Service” on page 2

“**acct\_get\_usage** Kernel Service” on page 3

### Related information:

acctctl Command

## add\_domain\_af Kernel Service

### Purpose

Adds an address family to the Address Family domain switch table.

### Syntax

```
#include <sys/types.h> #include <sys/errno.h> #include <sys/domain.h> int add_domain_af ( domain)
struct domain *domain;
```

### Parameter

Item	Description
<i>domain</i>	Specifies the domain of the address family.

### Description

The `add_domain_af` kernel service adds an address family domain to the Address Family domain switch table.

### Execution Environment

The `add_domain_af` kernel service can be called from either the process or interrupt environment.

### Return Values

Item	Description
0	Indicates that the address family was successfully added.
EEXIST	Indicates that the address family was already added.
EINVAL	Indicates that the address family number to be added is out of range.

### Example

To add an address family to the Address Family domain switch table, invoke the `add_domain_af` kernel service as follows:

```
add_domain_af(&inetdomain);
```

In this example, the family to be added is `inetdomain`.

#### Related reference:

“`del_domain_af` Kernel Service” on page 64

#### Related information:

Network Kernel Services

## add\_input\_type Kernel Service

### Purpose

Adds a new input type to the Network Input table.

### Syntax

```
#include <sys/types.h> #include <sys/errno.h> #include <net/if.h> #include <net/netisr.h> int
add_input_type (type, service_level, isr, ifq, af) u_short type; u_short service_level; int (* isr) (); struct
ifqueue * ifq; u_short af;
```



## Parameters

Item	Description
<i>type</i>	Specifies which type of protocol a packet contains. A value of x'FFFF' indicates that this input type is a wildcard type and matches all input packets.
<i>service_level</i>	Determines the processing level at which the protocol input handler is called. If the <i>service_level</i> parameter is set to <b>NET_OFF_LEVEL</b> , the input handler specified by the <i>isr</i> parameter is called directly. Setting the <i>service_level</i> parameter to <b>NET_KPROC</b> schedules a network dispatcher. This dispatcher calls the subroutine identified by the <i>isr</i> parameter.
<i>isr</i>	Identifies the routine that serves as the input handler for an input packet type.
<i>ifq</i>	Specifies an input queue for holding input buffers. If this parameter has a non-null value, an input buffer ( <b>mbuf</b> ) is enqueued. The <i>ifq</i> parameter must be specified if the processing level specified by the <i>service_level</i> parameter is <b>NET_KPROC</b> . Specifying null for this parameter generates a call to the input handler specified by the <i>isr</i> parameter, as in the following:
<i>af</i>	Specifies the address family of the calling protocol. The <i>af</i> parameter must be specified if the <i>ifq</i> parameter is not a null character. This parameter must be greater than or equal to 0 and less than <b>NETISR_MAX</b> . Refer to <b>netisr.h</b> for the range of values of <i>af</i> that are already in use. Also, other kernel extensions that are not AIX and that use network ISRs currently running on the system can make use of additional values not mentioned in <b>netisr.h</b> .  (*isr)(CommonPortion,Buffer);  In this example, CommonPortion points to the network common portion (the <b>arpcom</b> structure) of a network interface and Buffer is a pointer to a buffer ( <b>mbuf</b> ) containing an input packet.

## Description

To enable the reception of packets, an address family calls the **add\_input\_type** kernel service to register a packet type in the Network Input table. Multiple packet types require multiple calls to *Kernel Extensions and Device Support Programming Concepts* the **add\_input\_type** kernel service.

## Execution Environment

The **add\_input\_type** kernel service can be called from either the process or interrupt environment.

## Return Values

Item	Description
0	Indicates that the type was successfully added.
EEXIST	Indicates that the type was previously added to the Network Input table.
ENOSPC	Indicates that no free slots are left in the table.
EINVAL	Indicates that an error occurred in the input parameters.

## Examples

1. To register an Internet packet type (**TYPE\_IP**), invoke the **add\_input\_type** service as follows:

```
add_input_type(TYPE_IP, NET_KPROC, ipintr, &ipintrq, AF_INET);
```

This packet is processed through the network kproc. The input handler is ipintr. The input queue is ipintrq.

2. To specify the input handler for ARP packets, invoke the **add\_input\_type** service as follows:

```
add_input_type(TYPE_ARP, NET_OFF_LEVEL, arpinput, NULL, NULL);
```

Packets are not queued and the arpinput subroutine is called directly.

### Related reference:

“del\_input\_type Kernel Service” on page 65

“find\_input\_type Kernel Service” on page 143

### Related information:

## add\_netisr Kernel Service

### Purpose

Adds a network software interrupt service to the Network Interrupt table.

### Syntax

```
#include <sys/types.h> #include <sys/errno.h> #include <net/netisr.h> int add_netisr ( soft_intr_level, service_level, isr) u_short soft_intr_level; u_short service_level; int (*isr) ();
```

### Parameters

Item	Description
<i>soft_intr_level</i>	Specifies the software interrupt level to add. This parameter must be greater than or equal to 0 and less than <b>NETISR_MAX</b> . Refer to <b>netisr.h</b> for the range of values of <i>soft_intr_level</i> that are already in use. Also, other kernel extensions that are not AIX and that use network ISRs currently running on the system can make use of additional values not mentioned in <b>netisr.h</b> .
<i>service_level</i>	Specifies the processing level of the network software interrupt.
<i>isr</i>	Specifies the interrupt service routine to add.

### Description

The **add\_netisr** kernel service adds the software-interrupt level specified by the *soft\_intr\_level* parameter to the Network Software Interrupt table.

The processing level of a network software interrupt is specified by the *service\_level* parameter. If the interrupt level specified by the *service\_level* parameter equals **NET\_KPROC**, a network interrupt scheduler calls the function specified by the *isr* parameter. If you set the *service\_level* parameter to **NET\_OFF\_LEVEL**, the **schednetisr** service calls the interrupt service routine directly.

### Execution Environment

The **add\_netisr** kernel service can be called from either the process or interrupt environment.

### Return Values

Item	Description
0	Indicates that the interrupt service routine was successfully added.
EEXIST	Indicates that the interrupt service routine was previously added to the table.
EINVAL	Indicates that the value specified for the <i>soft_intr_level</i> parameter is out of range or at a service level that is not valid.

### Related reference:

“del\_netisr Kernel Service” on page 66

### Related information:

Network Kernel Services

## add\_netopt Macro

### Purpose

Adds a network option structure to the list of network options.

## Syntax

```
#include <sys/types.h> #include <sys/errno.h> #include <net/netopt.h> add_netopt (  
option_name_symbol, print_format) option_name_symbol; char *print_format;
```

## Parameters

Item	Description
<i>option_name_symbol</i>	Specifies the symbol name used to construct the <b>netopt</b> structure and default names.
<i>print_format</i>	Specifies the string representing the print format for the network option.

## Description

The **add\_netopt** macro adds a network option to the linked list of network options. The **no** command can then be used to show or alter the variable's value.

The **add\_netopt** macro has no return values.

## Execution Environment

The **add\_netopt** macro can be called from either the process or interrupt environment.

### Related reference:

“del\_netopt Macro” on page 67

### Related information:

no Command

Network Kernel Services

## as\_att64 Kernel Service

### Purpose

Allocates and maps a specified region in the current user address space.

## Syntax

```
#include <sys/types.h> #include <sys/errno.h> #include <sys/vmuser.h> #include <sys/adspace.h>  
unsigned long long as_att64 (vmhandle, offset) vmhandle_t vmhandle; int offset;
```

## Parameters

Item	Description
<i>vmhandle</i>	Describes the virtual memory object being made addressable in the address space.
<i>offset</i>	Specifies the offset in the virtual memory object. The upper 4-bits of this offset are ignored.

## Description

<b>Item</b>	<b>Description</b>
The <b>as_att64</b> kernel service:	<p>Selects an unallocated region within the current user address space.</p> <p>Allocates the region.</p> <p>Maps the virtual memory object selected by the <code>vmhandle</code> parameter with the access permission specified in the handle.</p> <p>Constructs the address of the offset specified by the <code>offset</code> parameter within the user-address space.</p>

The **as\_att64** kernel service assumes an address space model of fixed-size virtual memory objects.

This service will operate correctly for both 32-bit and 64-bit user address spaces. It will also work for kernel processes (*kprocs*).

**Note:** This service only operates on the current process's address space. It is not allowed to operate on another address space.

## Execution Environment

The **as\_att64** kernel service can be called from the process environment only.

## Return Values

On successful completion, this service returns the base address plus the input offset (`offset`) into the allocated region.

<b>Item</b>	<b>Description</b>
NULL	An error occurred and <code>errno</code> indicates the cause:
EINVAL	Address specified is out of range, or
ENOMEM	Could not allocate due to insufficient resources.

## Related reference:

“`as_seth64` Kernel Service” on page 21

“`as_geth64` Kernel Service” on page 14

“`as_getsrval64` Kernel Service” on page 15

## **as\_det64** Kernel Service

### Purpose

Unmaps and deallocates a region in the current user address space that was mapped with the **as\_att64** kernel service.

### Syntax

```
#include <sys/errno.h> #include <sys/adspace.h> int as_det64 (addr64) unsigned long long addr64;
```

### Parameters

Item	Description
<code>addr64</code>	Specifies an effective address within the region to be deallocated.

## Description

The `as_det64` kernel service unmaps the virtual memory object from the region containing the specified effective address (specified by the `addr64` parameter).

The `as_det64` kernel service assumes an address space model of fixed-size virtual memory objects.

This service should not be used to deallocate a base kernel region, process text, process private or an unallocated region. An `EINVAL` return code will result.

This service will operate correctly for both 32-bit and 64-bit user address spaces. It will also work for kernel processes (*kprocs*).

**Note:** This service only operates on the current process's address space. It is not allowed to operate on another address space.

## Execution Environment

The `as_det64` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	The region was successfully unmapped and deallocated.
<code>EINVAL</code>	An attempt was made to deallocate a region that should not have been deallocated (that is, a base kernel region, process text region, process private region, or unallocated region).
<code>EINVAL</code>	Input address out of range.

### Related reference:

“`as_seth64` Kernel Service” on page 21

“`as_geth64` Kernel Service” on page 14

“`as_getsrval64` Kernel Service” on page 15

## `as_geth` Kernel Service

### Purpose

Obtains a handle to the virtual memory object for the specified address given in the specified address space.

### Syntax

```
#include <sys/types.h> #include <sys/errno.h> #include <sys/vmuser.h> #include <sys/adspace.h>
vmhandle_t as_geth (Adspacep, Addr) adspace_t *Adspacep; caddr_t Addr;
```

### Parameters

Item	Description
<i>Adspacep</i>	Points to the address space structure to obtain the virtual memory object handle from.
<i>Addr</i>	Specifies the virtual memory address that should be used to determine the virtual memory object handle for the specified address space.

## Description

The **as\_geth** kernel service is used to obtain a handle to the virtual memory object corresponding to a virtual memory address in a particular address space. This handle can then be used with the **vm\_att** kernel service to make the object addressable in another address space.

This service can also be called from the interrupt environment.

## Execution Environment

The **as\_geth** kernel service can be called from the process environment only.

## Return Values

The **as\_geth** kernel service always succeeds and returns the appropriate handle.

### Related reference:

“vm\_att Kernel Service” on page 536

## as\_geth64 Kernel Service

### Purpose

Obtains a handle to the virtual memory object for the specified address.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
#include <sys/adspace.h>
```

```
vmhandle_t as_geth64 (addr64)
unsigned long long addr64;
```

### Parameter

Item	Description
<b>addr64</b>	Specifies the virtual memory address for which the corresponding handle should be returned.

## Description

The **as\_geth64** kernel service is used to obtain a handle to the virtual memory object corresponding to the input address (*addr64*). This handle can then be used with the **as\_att64** or **vm\_att** kernel service to make the object addressable at a different location.

After the last use of the handle and after it is detached accordingly, the **as\_puth64** kernel service must be used to indicate this fact. Failure to call the **as\_puth64** service may result in resources being permanently unavailable for re-use.

If the handle returned refers to a virtual memory segment, then that segment is protected from deletion until the **as\_puth64** kernel service is called.

If, for some reason, it is known that the virtual memory object cannot be deleted, then the **as\_getsrval64** kernel service may be used instead of the **as\_geth64** service.

The **as\_geth64** kernel service assumes an address space model of fixed-size virtual memory objects.

This service will operate correctly for both 32-bit and 64-bit user address spaces. It will also work for kernel processes (*kprocs*).

**Note:** This service only operates on the current process's address space. It is not allowed to operate on another address space.

## Execution Environment

The **as\_geth64** kernel service can be called from the process environment only.

## Return Values

On successful completion, this routine returns the appropriate handle.

On error, this routine returns the value `INVLSID` defined in **sys/seg.h**. This is caused by an address out of range.

Errors include: Input address out of range.

### Related reference:

“as\_seth64 Kernel Service” on page 21

“as\_getsrval64 Kernel Service”

“as\_puth64 Kernel Service” on page 20

## as\_getsrval64 Kernel Service

### Purpose

Obtains a handle to the virtual memory object for the specified address.

### Syntax

```
#include <sys/types.h> #include <sys/errno.h> #include <sys/vmuser.h> #include <sys/adspace.h>
vmhandle_t as_getsrval64 (addr64) unsigned long long addr64;
```

### Parameters

Item	Description
<i>addr64</i>	Specifies the virtual memory address for which the corresponding handle should be returned.

### Description

The **as\_getsrval64** kernel service is used to obtain a handle to the virtual memory object corresponding to the input address(*addr64*). This handle can then be used with the **as\_att64** or **vm\_att** kernel services to make the object addressable at a different location.

This service should only be used when it is known that the virtual memory object cannot be deleted, otherwise the **as\_geth64** kernel service must be used.

The **as\_puth64** kernel service must not be called for handles returned by the **as\_getsrval64** kernel service.

The **as\_getsrval64** kernel service assumes an address space model of fixed-size virtual memory objects.

This service will operate correctly for both 32-bit and 64-bit user address spaces. It will also work for kernel processes (*kprocs*).

**Note:** This service only operates on the current process's address space. It is not allowed to operate on another address space.

## Execution Environment

The `as_getsrval64` kernel service can be called from the process environment only when the current user address space is 64-bits. If the current user address space is 32-bits, or is a *kproc*, then `as_getsrval64` may be called from an interrupt environment.

## Return Values

On successful completion this routine returns the appropriate handle.

On error, this routine returns the value `INVLSID` defined in `sys/seg.h`. This is caused by an address out of range.

Errors include: Input address out of range.

### Related reference:

“`as_geth64` Kernel Service” on page 14

“`as_puth64` Kernel Service” on page 20

“`as_seth64` Kernel Service” on page 21

## as\_lw\_att64 Kernel Service

### Purpose

Allocates and maps a specified region in the current user address space. Part of the lightweight kernel service subsystem, which must be initialized with the `as_lw_pool_init` kernel service before it can be used.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sysvmuser.h>
#include <sys/adspace.h>
#include <sys/mem.h>
```

```
int as_lw_att64 (dp, offset, length, addr)
xmem* dp;
size_t offset;
size_t length;
ptr64* addr;
```

### Parameters

Item	Description
<i>dp</i>	Pointer to a cross memory descriptor that describes the virtual memory object that is being made addressable in the address space.
<i>offset</i>	Specifies the byte offset in the virtual memory object.
<i>length</i>	Specifies the number of bytes to map in the virtual memory object.
<i>addr</i>	Pointer to the location where the address will be returned.



## Description

The `as_lw_att64` kernel service does the following:

- Allocates a region from the process' address space for the mapping.
- Maps the virtual memory object selected by the `dp` parameter.
- Constructs the address of the offset specified by the `offset` parameter within the user-address space.

**Note:** The `as_lw_att64` kernel service should be used with caution. Be sure to read the documentation for this and the other lightweight services (`as_lw_det64` and `as_lw_pool_init`) carefully before doing so. There is a risk of illegal data access and cross-process data corruption if these services are not used correctly.

In order to use this service, the cross memory descriptor pointed to by the `dp` parameter must be initialized by using the `xmattach` kernel service with the `LW_XMATTACH` flag set. The `lw_pool_init` kernel service must also have been successfully called by the current process.

The service will map an area length *bytes* long into the caller's address space from the memory represented by the descriptor, starting at the number of bytes specified in the `offset` parameter. It is illegal for any thread other than the caller of this service to address the attached region.

This service will operate correctly only in 64-bit user address spaces. It will not work for kernel processes (kprocs).

**Note:** This service only operates on the current process's address space. It is not allowed to operate on another address space.

## Execution Environment

The `as_lw_att64` kernel service can be called from the process environment only.

## Return Values

On successful completion, this service sets the value of `addr` to the address of the allocated region and returns 0.

Item	Description
NULL	An error occurred and <code>errno</code> indicates the cause.
EINVAL	Cross memory descriptor is in an invalid state, length is zero or offset plus length goes past the end of the virtual memory object.
ENODEV	The <code>as_lw_pool_init</code> kernel service has not been called to initialize the pool settings for this process.
ENOSYS	Called by a 32-bit process.
ENOSPC	Resources allocated to do lightweight services for this thread expended. Either the region to be attached is too large (the <code>as_lw_pool_init</code> kernel service was called with too small a <code>pool_size</code> ) or there are outstanding attaches which need to release their lightweight resources using the <code>as_lw_det64</code> kernel service before this attach can be completed.
EIO	Indicates a failure of the lightweight subsystem, process should discontinue use of lightweight kernel services.
EPERM	Called by a user thread that is not 1:1 with a kernel thread.
ENOMEM	Could not allocate system resources for lightweight services for this thread.

## Implementation Specifics

The `as_lw_att64` kernel service is part of Base Operating System (BOS) Runtime.

**Related reference:**

“`as_lw_det64` Kernel Service” on page 18

“as\_lw\_pool\_init Kernel Service” on page 19

## as\_lw\_det64 Kernel Service

### Purpose

Unmaps and deallocates a region in the current user address space that was mapped using the `as_lw_att64` kernel service.

### Syntax

```
#include <sys/errno.h>
#include <sys/adspc.h>
#include <sys/xmem.h>
int as_lw_det64 (dp, addr, length)
xmem* dp;
void* addr;
size_t length;
```

### Parameters

Item	Description
<i>dp</i>	The cross memory descriptor describing the attached virtual memory.
<i>addr</i>	Specifies the first effective address of the region to be deallocated.
<i>length</i>	Specifies the length of the region to be deallocated.

### Description

**Note:** The `as_lw_det64` kernel service should be used with caution. Read the documentation for this and the other lightweight services (`as_lw_att64` and `as_lw_pool_init`) carefully before doing so. There is a risk that illegal data accesses will be allowed if these services are not used correctly.

The `as_lw_det64` kernel service unmaps the virtual memory from the region starting at the specified effective address, which is specified by the *addr* parameter. This service (and only this service) must be used to unmap regions mapped by the `as_lw_att64` kernel service. It must be called by the same thread that called the `as_lw_att64` kernel service. The *addr* parameter must be the value returned by the `as_lw_att64` kernel service, and the *dp* parameter and the *length* parameter must be the same *dp* and *length* passed to it. The `xmdetach` kernel service must not be called to release the *dp* parameter until any outstanding attaches of the *dp* parameter using the `as_lw_att64` kernel service have been detached using the `as_lw_det64` kernel service.

The `as_lw_det64` kernel service cannot be used to detach a region not mapped by the `as_lw_att64` kernel service.

The `as_lw_det64` kernel service will operate correctly only for 64-bit user address spaces. It will not work for kernel processes (kprocs).

**Note:** This service only operates on the current process's address space. It is not allowed to operate on another address space.

### Execution Environment

The `as_lw_det64` kernel service can be called from the process environment only.

### Return Values

Item	Description
0	The region was successfully unmapped and deallocated.
EINVAL	An attempt was made to deallocate a region that should not have been deallocated.
ENOSYS	The service was called by a 32-bit process.
ENOMEM	No lightweight resources allocated to this thread.
EIO	Indicates a failure of the lightweight subsystem, process should discontinue use of lightweight kernel services.
EPERM	Called by a user thread that is not 1:1 with a kernel thread.

## Implementation Specifics

The `as_lw_det64` kernel service is part of Base Operating System (BOS) Runtime.

### Related reference:

“`as_lw_att64` Kernel Service” on page 16

“`as_lw_pool_init` Kernel Service”

## as\_lw\_pool\_init Kernel Service

### Purpose

Initializes lightweight attach and detach subsystem for the current process with the given settings.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
#include <sys/adspace.h>
```

```
int as_lw_pool_init (pool_size, flags)
size_t pool_size;
uint flags;
```

### Parameters

Item	Description
<i>pool_size</i>	Specifies the maximum number of bytes that can be attached by lightweight services at one time by each thread of this process.
<i>flags</i>	Specifies flag options for this kernel service. Valid values are 0 and LW_DEBUG.

### Description

**Note:** The `as_lw_pool_init` kernel service should be used with caution. Read the documentation for this and the other lightweight services (`as_lw_att64` and `as_lw_det64`) carefully before doing so. There is a risk that illegal data accesses will be allowed if these services are not used correctly.

The `as_lw_pool_init` kernel service initializes the lightweight pool size and flag settings for the current process. Once it has been called, these settings are fixed and cannot be changed for the process.

If LW\_DEBUG is set in the *flags* parameter, the risk of illegal data access will be removed from calls to the `as_lw_att64` kernel service and the `as_lw_det64` kernel service. This setting allows users to debug problems that are caused by incorrect use of these services.

Processes that have called the `as_lw_pool_init` kernel service can use the other lightweight kernel services (`as_lw_att64` and `as_lw_det64`) to attach and detach virtual memory regions represented by a cross memory descriptor. These kernel services are used on a per-thread basis, that is if one thread uses the

**as\_lw\_att64** kernel service to attach virtual memory to a region of its address space, that region cannot be addressed by any other thread, and it must be detached by the same thread by using the **as\_lw\_det64** kernel service.

This service will operate correctly only for 64-bit user address spaces. It will not work for kernel processes (kprocs).

## Execution Environment

The **as\_lw\_pool\_init** kernel service can be called from a 64-bit process environment only.

## Return Values

On successful completion, this service returns 0.

Item	Description
ENOSYS	The service was called by a 32-bit process.
EEXIST	The <b>as_lw_pool_init</b> kernel service has already been successfully completed for this process.
EINVAL	Invalid flag settings or the <i>pool_size</i> parameter is 0.
EPERM	Called by a user thread that is not 1:1 with a kernel thread.

## Implementation Specifics

The **as\_lw\_pool\_init** kernel service is part of Base Operating System (BOS) Runtime.

### Related reference:

“**as\_lw\_att64** Kernel Service” on page 16

“**as\_lw\_det64** Kernel Service” on page 18

## as\_puth64 Kernel Service

### Purpose

Indicates that no more references will be made to a virtual memory object obtained using the **as\_geth64** kernel service.

### Syntax

```
#include <sys/types.h> #include <sys/errno.h> #include <sys/vmuser.h> #include <sys/adspace.h> int  
as_puth64 ( addr64, vmhandle ) unsigned long long addr64; vmhandle_t vmhandle;
```

### Parameters

Item	Description
<i>addr64</i>	Specifies the virtual memory address that the virtual memory object handle was obtained from. This must be the same address that was given to the <b>as_geth64</b> kernel service previously.
<i>vmhandle</i>	Describes the virtual memory object that will no longer be referenced. This handle must have been returned by the <b>as_geth64</b> kernel service.

### Description

The **as\_puth64** kernel service is used to indicate that no more references will be made to the virtual memory object returned by a call to the **as\_geth64** kernel service. The virtual memory object must be detached from the address space already, using either **as\_det64** or **vm\_det** service.

Failure to call the **as\_puth64** kernel service may result in resources being permanently unavailable for re-use.

If, for some reason, it is known that the virtual memory object cannot be deleted, the `as_getsrval64` kernel service may be used instead of the `as_geth64` kernel service. This kernel service does not require that the `as_puth64` kernel service be used.

The `as_puth64` kernel service assumes an address space model of fixed-size virtual memory objects.

This service will operate correctly for both 32-bit and 64-bit user address spaces. It will also work for kernel processes (*kprocs*).

**Note:** This service only operates on the current process's address space. It is not allowed to operate on another address space.

## Execution Environment

The `as_puth64` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Successful completion.
EINVAL	Input address out of range.

### Related reference:

“`as_getsrval64` Kernel Service” on page 15

“`as_geth64` Kernel Service” on page 14

“`as_seth64` Kernel Service”

## `as_seth64` Kernel Service

### Purpose

Maps a specified region for the specified virtual memory object.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
#include <sys/adspace.h>
```

```
int as_seth64 (addr64, vmhandle)
unsigned long long addr64;
vmhandle_t vmhandle;
```

### Parameters

Item	Description
<i>addr64</i>	The region covering this input virtual memory address will be mapped.
<i>vmhandle</i>	Describes the virtual memory object being made addressable within a region of the address space.

### Description

The `as_seth64` kernel service maps the region covering the input `addr64` parameter. Any virtual memory object previously mapped within this region is unmapped.

The virtual memory object specified with the `vmhandle` parameter is then mapped with the access permission specified in the handle.

The **as\_seth64** kernel service should only be used when it is necessary to map a virtual memory object at a fixed address. The **as\_att64** kernel service should be used when it is not absolutely necessary to map the virtual memory object at a fixed address.

The **as\_seth64** kernel service assumes an address space model of fixed-size virtual memory objects.

This service will operate correctly for both 32-bit and 64-bit user address spaces. It will also work for kernel processes (*kprocs*).

**Note:** This service only operates on the current process's address space. It is not allowed to operate on another address space.

## Execution Environment

The **as\_seth64** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Successful completion.
EINVAL	Input address out of range.

### Related reference:

“**as\_det64** Kernel Service” on page 12

“**as\_geth64** Kernel Service” on page 14

“**as\_puth64** Kernel Service” on page 20

## attach Device Queue Management Routine

### Purpose

Provides a means for performing device-specific processing when the **attchq** kernel service is called.

### Syntax

```
#include <sys/types.h> #include <sys/errno.h> #include <sys/deviceq.h> int attach ( dev_parms, path_id)
caddr_t dev_parms; cba_id path_id;
```

### Parameters

Item	Description
<i>dev_parms</i>	Passed to the <b>creatd</b> kernel service when the <b>attach</b> routine is defined.
<i>path_id</i>	Specifies the path identifier for the queue being attached to.

### Description

The **attach** routine is part of the Device Queue Management kernel extension. Each device queue can have an **attach** routine. This routine is optional and must be specified when the **creatd** kernel service defines the device queue. The **attchq** service calls the **attach** routine each time a new path is created to the owning device queue. The processing performed by this routine is dependent on the server function.

The **attach** routine executes under the process under which the **attchq** kernel service is called. The kernel does not serialize the execution of this service with the execution of any other server routines.

## Execution Environment

The `attach-device` routine can be called from the process environment only.

## Return Values

Item	Description
<code>RC_GOOD</code>	Indicates a successful completion.
<code>RC_NONE</code>	Indicates that resources such as pinned memory are unavailable.
<code>RC_MAX</code>	Indicates that the server already has the maximum number of users that it supports.
Greater than or equal to <code>RC_DEVICE</code>	Indicates device-specific errors.

## audit\_svcbcopy Kernel Service

### Purpose

Appends event information to the current audit event buffer.

### Syntax

```
#include <sys/types.h> #include <sys/errno.h> int audit_svcbcopy ( buf, len) char *buf; int len;
```

### Parameters

Item	Description
<i>buf</i>	Specifies the information to append to the current audit event record buffer.
<i>len</i>	Specifies the number of bytes in the buffer.

### Description

The `audit_svcbcopy` kernel service appends the specified buffer to the event-specific information for the current switched virtual circuit (SVC). System calls should initialize auditing with the `audit_svcstart` kernel service, which creates a record buffer for the named event.

The `audit_svcbcopy` kernel service can then be used to add additional information to that buffer. This information usually consists of system call parameters passed by reference.

If auditing is enabled, the information is written by the `audit_svcfinis` kernel service after the record buffer is complete.

## Execution Environment

The `audit_svcbcopy` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
ENOMEM	Indicates that the kernel service is unable to allocate space for the new buffer.

**Related reference:**

“audit\_svcfinis Kernel Service”  
“audit\_svcstart Kernel Service”

**Related information:**

Security Kernel Services

## audit\_svcfinis Kernel Service

### Purpose

Writes an audit record for a kernel service.

### Syntax

```
#include <sys/types.h> #include <sys/errno.h> #include <sys/audit.h> int audit_svcfinis ( )
```

### Description

The **audit\_svcfinis** kernel service completes an audit record begun earlier by the **audit\_svcstart** kernel service and writes it to the kernel audit logger. Any space allocated for the record and associated buffers is freed.

If the system call terminates without calling the **audit\_svcfinis** service, the switched virtual circuit (SVC) handler exit routine writes the records. This exit routine calls the **audit\_svcfinis** kernel service to complete the records.

### Execution Environment

The **audit\_svcfinis** kernel service can be called from the process environment only.

### Return Values

The **audit\_svcfinis** kernel service always returns a value of 0.

**Related reference:**

“audit\_svcbcopy Kernel Service” on page 23  
“audit\_svcstart Kernel Service”

**Related information:**

Security Kernel Services

## audit\_svcstart Kernel Service

### Purpose

Initiates an audit record for a system call.

### Syntax

```
#include <sys/types.h> #include <sys/errno.h> #include <sys/audit.h> int audit_svcstart (eventnam  

, eventnum, numargs, arg1, arg2, ...) char * eventnam; int * eventnum; int numargs; int arg1; int arg2; ...
```



## Parameters

Item	Description
<i>eventnam</i>	Specifies the name of the event. In the current implementation, event names must be less than 17 characters, including the trailing null character. Longer names are truncated.
<i>eventnum</i>	Specifies the number of the event. This is an internal table index meaningful only to the kernel audit logger. The system call should initialize this parameter to 0. The first time the <b>audit_svcstart</b> kernel service is called, this parameter is set to the actual table index. The system call should not reset the parameter. The parameter should be declared a static.
<i>numargs</i>	Specifies the number of parameters to be included in the buffer for this record. These parameters are normally zero or more of the system call parameters, although this is not a requirement.
<i>arg1, arg2, ...</i>	Specifies the parameters to be included in the buffer.

## Description

The **audit\_svcstart** kernel service initiates auditing for a system call event. It dynamically allocates a buffer to contain event information. The arguments to the system call (which should be specified as parameters to this kernel service) are automatically added to the buffer, as is the internal number of the event. You can use the **audit\_svbcopy** service to add additional information that cannot be passed by value.

The system call commits this record with the **audit\_svcfinis** kernel service. The system call should call the **audit\_svcfinis** kernel service before calling another system call.

## Execution Environment

The **audit\_svcstart** kernel service can be called from the process environment only.

## Return Values

Item	Description
<b>Nonzero</b>	Indicates that auditing is on for this routine.
<b>0</b>	Indicates that auditing is off for this routine.

## Example

```
svccrash(int x, int y, int z)
{
    static int eventnum;
    if (audit_svcstart("crashed", &eventnum, 2, x, y))
    {
        audit_svcfinis();
    }
    body of svccrash
}
```

The preceding example allocates an audit event record buffer for the crashed event and copies the first and second arguments into it. The third argument is unnecessary and not copied.

### Related reference:

“audit\_svbcopy Kernel Service” on page 23

“audit\_svcfinis Kernel Service” on page 24

### Related information:

Security Kernel Services

## b

The following kernel services begin with the with the letter b.

## **bawrite Kernel Service**

### **Purpose**

Writes the specified buffer data without waiting for I/O to complete.

### **Syntax**

```
#include <sys/types.h> #include <sys/errno.h> #include <sys/buf.h> int bawrite ( bp) struct buf *bp;
```

### **Parameter**

Item	Description
------	-------------

<i>bp</i>	Specifies the address of the buffer structure.
-----------	--

On a platform that supports storage keys, the passed in *bp* parameter must be in the **KKEY\_PUBLIC** or **KKEY\_BLOCK\_DEV** protection domain.

### **Description**

The **bawrite** kernel service sets the asynchronous flag in the specified buffer and calls the **bwrite** kernel service to write the buffer.

### **Execution Environment**

The **bawrite** kernel service can be called from the process environment only.

### **Return Values**

Item	Description
------	-------------

0	Indicates successful completion.
---	----------------------------------

ERRNO	Returns an error number from the <code>/usr/include/sys/errno.h</code> file on error.
-------	---

### **Related reference:**

“bwrite Kernel Service” on page 36

### **Related information:**

Block I/O buffer cache kernel services overview

I/O Kernel Services

## **bdwrite Kernel Service**

### **Purpose**

Releases the specified buffer after marking it for delayed write.

### **Syntax**

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/buf.h>
```

```
void bdwrite ( bp)
struct buf *bp;
```

### **Parameter**

Item	Description
------	-------------

<i>bp</i>	Specifies the address of the buffer structure for the buffer to be written.
-----------	---

On a platform that supports storage keys, the passed in *bp* parameter must be in the **KKEY\_PUBLIC** or **KKEY\_BLOCK\_DEV** protection domain.

## Description

The **bdwrite** kernel service marks the specified buffer so that the block is written to the device when the buffer is stolen. The **bdwrite** service marks the specified buffer as delayed write and then releases it (that is, puts the buffer on the free list). When this buffer is reassigned or reclaimed, it is written to the device.

## Execution Environment

The **bdwrite** kernel service can be called from the process environment only.

## Return Values

The **bdwrite** kernel service has no return values.

### Related reference:

“brelse Kernel Service” on page 32

### Related information:

Block I/O Buffer Cache Kernel Services: Overview

I/O Kernel Services

## bflush Kernel Service

### Purpose

Flushes all write-behind blocks on the specified device from the buffer cache.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/buf.h>
```

```
void bflush ( dev)
dev_t dev;
```

### Parameter

Item	Description
------	-------------

<i>dev</i>	Specifies which device to flush. A value of <b>NODEVICE</b> flushes all devices.
------------	--

## Description

The **bflush** kernel service runs the free list of buffers. It notes as busy or writing any dirty buffer whose block is on the specified device. When a value of **NODEVICE** is specified, the **bflush** service flushes all write-behind blocks for all devices. The **bflush** service has no return values.

## Execution Environment

The **bflush** kernel service can be called from the process environment only.

### Related reference:

“bwrite Kernel Service” on page 36

## Related information:

Block I/O Buffer Cache Kernel Services: Overview

I/O Kernel Services

## bindprocessor Kernel Service

### Purpose

Binds or unbinds kernel threads to a processor.

### Syntax

```
#include <sys/processor.h>
```

```
int bindprocessor ( What, Who, Where )
```

```
int What;
```

```
int Who;
```

```
cpu_t Where;
```

### Parameters

Item	Description
------	-------------

<i>What</i>	Specifies whether a process or a kernel thread is being bound to a processor. The <i>What</i> parameter can take one of the following values:
-------------	---

**BINDPROCESS**

A process is being bound to a processor.

**BINDTHREAD**

A kernel thread is being bound to a processor.

<i>Who</i>	Indicates a process or kernel thread identifier, as appropriate for the <i>What</i> parameter, specifying the process or kernel thread which is to be bound to a processor.
------------	---

<i>Where</i>	If the <i>Where</i> parameter is in the range 0- <i>n</i> (where <i>n</i> is the number of online processors in the system), it represents a bind CPU identifier to which the process or kernel thread is to be bound. Otherwise, it represents a processor class, from which a processor will be selected. A value of <b>PROCESSOR_CLASS_ANY</b> unbinds the specified process or kernel thread, which will then be able to run on any processor.
--------------	--

### Description

The **bindprocessor** kernel service binds a single kernel thread, or all kernel threads in a process, to a processor, forcing the bound threads to be scheduled to run on that processor only. It is important to understand that a process itself is not bound, but rather its kernel threads are bound. Once kernel threads are bound, they are always scheduled to run on the chosen processor, unless they are later unbound. When a new thread is created using the **thread\_create** kernel service, it has the same bind properties as its creator.

Programs that use processor bindings must be aware of Dynamic Logical Partitioning (DLPAR).

### Return Values

On successful completion, the **bindprocessor** kernel service returns 0. Otherwise, a value of -1 is returned and the error code can be checked by calling the **getuerror** kernel service.

### Error Codes

The **bindprocessor** kernel service is unsuccessful if one of the following is true:

Item	Description
EINVAL	The <i>What</i> parameter is invalid, or the <i>Where</i> parameter indicates an invalid processor number or a processor class which is not currently available.
ESRCH	The specified process or thread does not exist.
EPERM	The caller does not have root user authority, and the <i>Who</i> parameter specifies either a process, or a thread belonging to a process, having a real or effective user ID different from that of the calling process.

## Execution Environment

The **bindprocessor** kernel service can be called from the process environment only.

### Related information:

bindprocessor command  
 fork subroutine  
 sysconf subroutine  
 Dynamic Logical Partitioning

## binval Kernel Service

### Purpose

Makes nonreclaimable all blocks in the buffer cache of a specified device.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/buf.h>
```

```
void binval ( dev)
dev_t dev;
```

### Parameter

Item	Description
<i>dev</i>	Specifies the device to be purged.

### Description

The **binval** kernel service makes nonreclaimable all blocks in the buffer cache of a specified device. Before removing the device from the system, use the **binval** service to remove the blocks.

All of blocks of the device to be removed need to be flushed before you call the **binval** service. Typically, these blocks are flushed after the last close of the device.

## Execution Environment

The **binval** kernel service can be called from the process environment only.

### Return Values

The **binval** service has no return values.

### Related reference:

“bflush Kernel Service” on page 27

### Related information:

Block I/O Buffer Cache Kernel Services Overview

## blkflush Kernel Service

### Purpose

Flushes the specified block if it is in the buffer cache.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/buf.h>
```

```
int blkflush ( dev, blkno)
dev_t dev;
daddr_t blkno;
```

### Parameters

Item	Description
<i>dev</i>	Specifies the device containing the block to be flushed.
<i>blkno</i>	Specifies the block to be flushed.

### Description

The **blkflush** kernel service checks to see if the specified buffer is in the buffer cache. If the buffer is not in the cache, then the **blkflush** service returns a value of 0. If the buffer is in the cache, but is busy, the **blkflush** service calls the **e\_sleep** service to wait until the buffer is no longer in use. Upon waking, the **blkflush** service tries again to access the buffer.

If the buffer is in the cache and is not busy, but is dirty, then it is removed from the free list. The buffer is then marked as busy and synchronously written to the device. If the buffer is in the cache and is neither busy nor dirty (that is, the buffer is already clean and therefore does not need to be flushed), the **blkflush** service returns a value of 0.

### Execution Environment

The **blkflush** kernel service can be called from the process environment only.

### Return Values

Item	Description
1	Indicates that the block was successfully flushed.
0	Indicates that the block was not flushed. The specified buffer is either not in the buffer cache or is in the buffer cache but neither busy nor dirty.

### Related reference:

“bwrite Kernel Service” on page 36

### Related information:

Block I/O Buffer Cache Kernel Services: Overview

I/O Kernel Services

## bread Kernel Service

### Purpose

Reads the specified block data into a buffer.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/buf.h>
```

```
struct buf *bread ( dev, blkno)
dev_t dev;
daddr_t blkno;
```

## Parameters

Item	Description
<i>dev</i>	Specifies the device containing the block to be read.
<i>blkno</i>	Specifies the block to be read.

## Description

The **bread** kernel service assigns a buffer to the given block. If the specified block is already in the buffer cache, then the block buffer header is returned. Otherwise, a free buffer is assigned to the specified block and the data is read into the buffer. The **bread** service waits for I/O to complete to return the buffer header.

The buffer is allocated to the caller and marked as busy.

## Execution Environment

The **bread** kernel service can be called from the process environment only.

## Return Values

The **bread** service returns the address of the selected buffer's header. A nonzero value for **B\_ERROR** in the **b\_flags** field of the buffer's header (**buf** structure) indicates an error. If this occurs, the caller should release the buffer associated with the block using the **brlse** kernel service.

On a platform that supports storage keys, the buffer header is allocated from the storage protected by the **KKEY\_BLOCK\_DEV** kernel key.

### Related reference:

“iowait Kernel Service” on page 233

### Related information:

Block I/O Buffer Cache Kernel Services: Overview

I/O Kernel Services

## breada Kernel Service

### Purpose

Reads in the specified block and then starts I/O on the read-ahead block.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/buf.h>
```

```

struct buf *breada ( dev, blkno, rablkno)
dev_t   dev;
daddr_t blkno;
daddr_t rablkno;

```

## Parameters

Item	Description
<i>dev</i>	Specifies the device containing the block to be read.
<i>blkno</i>	Specifies the block to be read.
<i>rablkno</i>	Specifies the read-ahead block to be read.

## Description

The **breada** kernel service assigns a buffer to the given block. If the specified block is already in the buffer cache, then the **bread** service is called to:

- Obtain the block.
- Return the buffer header.

Otherwise, the **getblk** service is called to assign a free buffer to the specified block and to read the data into the buffer. The **breada** service waits for I/O to complete and then returns the buffer header.

I/O is also started on the specified read-ahead block if the free list is not empty and the block is not already in the cache. However, the **breada** service does not wait for I/O to complete on this read-ahead block.

## Execution Environment

The **breada** kernel service can be called from the process environment only.

## Return Values

The **breada** service returns the address of the selected buffer's header. A nonzero value for **B\_ERROR** in the **b\_flags** field of the buffer header (**buf** structure) indicates an error. If this occurs, the caller should release the buffer associated with the block using the **brlse** kernel service.

On a platform that supports storage keys, the buffer header is allocated from the storage protected by the **KKEY\_BLOCK\_DEV** kernel key.

### Related reference:

“bread Kernel Service” on page 30

“iowait Kernel Service” on page 233

### Related information:

Block I/O Buffer Cache Kernel Services: Overview

## brlse Kernel Service

### Purpose

Frees the specified buffer.

### Syntax

```

#include <sys/types.h>
#include <sys/errno.h>
#include <sys/buf.h>

```



```
void brelse ( bp)
struct buf *bp;
```

## Parameter

Item	Description
------	-------------

<i>bp</i>	Specifies the address of the <b>buf</b> structure to be freed.
-----------	--

On a platform that supports storage keys, the passed in *bp* parameter must be in the **KKEY\_PUBLIC** or **KKEY\_BLOCK\_DEV** protection domain.

## Description

The **brelse** kernel service frees the buffer to which the *bp* parameter points.

The **brelse** kernel service awakens any processes waiting for this buffer or for another free buffer. The buffer is then put on the list of available buffers. The buffer is also marked as not busy so that it can either be reclaimed or reallocated.

The **brelse** service has no return values.

## Execution Environment

The **brelse** kernel service can be called from either the process or interrupt environment.

### Related reference:

“geteblk Kernel Service” on page 186

“buf Structure” on page 615

### Related information:

I/O Kernel Services

## bsr\_alloc Kernel Service Purpose

Allocates a Barrier Synchronization Register (BSR) resource, and retrieves mapping information.

## Syntax

```
#include <sys/adspace.h>
```

```
int bsr_alloc (
    int bsr_bytes,
    struct io_map * bsr_map,
    int *bsr_stride,
    int *bsr_id)
```

## Parameters

Item	Description
<i>bsr_bytes</i>	Number of BSR bytes wanted.
<i>bsr_map</i>	Mapping information for the BSR facility
<i>bsr_stride</i>	Stride at which the BSR bytes repeat within the mapping
<i>bsr_id</i>	An opaque identifier for the allocated BSR resource

## Description

The **bsr\_alloc** service can be used to allocate and reserve all or a portion of the BSR facility. The requested number of BSR bytes to allocate is communicated through the *bsr\_bytes* parameter. The

requested number of bytes must correspond to a supported window size, as communicated by the *supported\_window\_mask* parameter of the **bsr\_query** service. If the requested number of bytes is available, the bytes are reserved and the I/O mapping information for accessing the allocated facility is written to the *bsr\_map* structure. In addition, the stride within the mapping that the allocated BSR bytes repeat is recorded in the *bsr\_stride* field. The *bsr\_id* field is written with a unique identifier to be used with the **bsr\_free** call.

If multiple granules or windows are to be used, they must be allocated with independent calls to **bsr\_alloc**. This is because I/O mappings for multiple granules might not be contiguous, and strides are only applicable within the granule.

The resulting *bsr\_map* information can then be used as input to **rmmmap\_create** for establishing addressability to the BSR resource within the current process address space.

## Execution Environment

The **bsr\_alloc** service can only be called from the process environment.

## Return Values

If successful, **bsr\_alloc** returns 0 and modifies the *bsr\_map* structure so that it contains the mapping information for the newly allocated resource, modifies the *bsr\_stride* field displays the stride on which the BSR bytes repeat within the mapping, and modifies the *bsr\_id* field so that it displays a unique identifier for the newly allocated BSR resource. If unsuccessful, one of the following values is returned:

Item	Description
ENODEV	The BSR facility does not exist.
EINVAL	Unsupported number of bytes requested.
EBUSY	Requested BSR bytes or mappable BSR windows are currently in use.

## Related reference:

“bsr\_free Kernel Service”

“bsr\_query Kernel Service” on page 35

“rmmmap\_create Kernel Service” on page 450

## bsr\_free Kernel Service

### Purpose

Frees a Barrier Synchronization Register (BSR) resource previously allocated with the **bsr\_alloc** kernel service.

### Syntax

```
#include <sys/adspace.h>
```

```
int bsr_free (  
    int bsr_id,
```

### Parameters

Item	Description
<i>bsr_id</i>	BSR resource identifier as returned in the <i>bsr_id</i> field of the <b>bsr_alloc</b> call.

## Description

The **bsr\_free** service releases a BSR allocation. The specific BSR resource being freed is identified by the unique identifier *bsr\_id* from the corresponding **bsr\_alloc** call.

It is the caller's responsibility to ensure that all prior attachments to the BSR resource, through **rmmap\_create** calls, have been detached with corresponding **rmmap\_remove** calls prior to freeing the BSR resource.

## Execution Environment

The **bsr\_free** service can only be called from the process environment.

## Return Values

Item	Description
0	A successful operation.
ENODEV	The BSR facility is not present.
EINVAL	BSR resource corresponding to <i>bsr_id</i> is invalid or not currently allocated.

### Related reference:

“bsr\_alloc Kernel Service” on page 33

“bsr\_query Kernel Service”

“rmmap\_remove Kernel Service” on page 454

## bsr\_query Kernel Service

### Purpose

Queries the existence of the Barrier Synchronization Register facility, and, if it exists, its size and allocation granule.

### Syntax

```
#include <sys/adspace.h>
```

```
int bsr_query (
    int *total_bytes,
    uint *supported_window_mask,
    int *free_bytes,
    uint *free_window_mask)
```

### Parameters

Item	Description
<i>total_bytes</i>	Total bytes of the BSR facility currently present within the system or logical partition
<i>supported_window_mask</i>	Bit mask representing supported power-of-2-sized windows that can be allocated
<i>free_bytes</i>	Number of BSR bytes currently available (not allocated)
<i>free_window_mask</i>	Bit mask representing available (not allocated) power-of-2-sized windows

## Description

The `bsr_query` service can be used to detect the presence and capabilities of the Barrier Synchronization Register (BSR) facility on a given system or logical partition. If the BSR facility is present on a system or within a logical partition, a value of 0 is returned, and the parameters, passed by reference, are written with the appropriate information.

The `total_bytes` field is written with the total number of BSR bytes currently present in the system or logical partition. The `supported_window_mask` field is written with a bitmask, where each bit set indicates the various power-of-2 window sizes that the `total_bytes` can be allocated and accessed. For example, a mask of 0x58 would indicate that windows of size 64 (0x40), 16 (0x10), and 8 (0x8) bytes were supported.

The `free_bytes` field is written with the number of BSR bytes within the system or logical partition that are currently unallocated. The `free_window_mask` field is written with a bitmask, where each bit set indicates the power-of-2 window sizes that are available for allocating and accessing the remaining `free_bytes`.

**Note:** Due to dynamic reconfiguration, the information returned by this query service might become stale.

## Execution Environment

The `bsr_query` service can only be called from the process environment.

## Return Values

Item	Description
0	The BSR facility exists and information is provided.
ENODEV	The BSR facility does not exist.

### Related reference:

“`bsr_alloc` Kernel Service” on page 33

“`bsr_free` Kernel Service” on page 34

## bwrite Kernel Service

### Purpose

Writes the specified buffer data.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/buf.h>
```

```
int bwrite ( bp)
struct buf *bp;
```

### Parameter

Item	Description
------	-------------

<i>bp</i>	Specifies the address of the buffer structure for the buffer to be written.
-----------	---

On a platform that supports storage keys, the passed in *bp* parameter must be in the **KKEY\_PUBLIC** or **KKEY\_BLOCK\_DEV** protection domain.

## Description

The **bwrite** kernel service writes the specified buffer data. If this is a synchronous request, the **bwrite** service waits for the I/O to complete.

"Block I/O Buffer Cache Kernel Services: Overview" in *Kernel Extensions and Device Support Programming Concepts* describes how the three buffer-cache write routines work.

## Execution Environment

The **bwrite** kernel service can be called from the process environment only.

## Return Values

Item	Description
------	-------------

0	Indicates a successful operation.
---	-----------------------------------

ERRNO	Returns an error number from the <code>/usr/include/sys/errno.h</code> file on error.
-------	---

### Related reference:

"brelse Kernel Service" on page 32

"iowait Kernel Service" on page 233

### Related information:

I/O Kernel Services

## C

The following kernel services begin with the with the letter c.

## cancel Device Queue Management Routine

### Purpose

Provides a means for cleaning up queue element-related resources when a pending queue element is eliminated from the queue.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/deviceq.h>
```

```
void cancel ( ptr)
struct req_qe *ptr;
```

### Parameter

Item	Description
<i>ptr</i>	Specifies the address of the queue element.

## Description

The kernel calls the **cancel** routine to clean up resources associated with a queue element. Each device queue can have a **cancel** routine. This routine is optional and must be specified when the device queue is created with the **creatq** service.

The **cancel** routine is called when a pending queue element is eliminated from the queue. This occurs when the path is destroyed or when the **cancelq** service is called. The device manager should unpin any data and detach any cross-memory descriptor.

Any operations started as a result of examining the queue with the **peekq** service must be stopped.

The **cancel** routine is also called when a queue is destroyed to get rid of any pending or active queue elements.

## Execution Environment

The **cancel-queue-element** routine can be called from the process environment only.

## cfgnadd Kernel Service

### Purpose

Registers a notification routine to be called when system-configurable variables are changed.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/sysconfig.h>
```

```
void cfgnadd
( cbp)
struct cfgncb *cbp;
```

### Parameter

Item	Description
<i>cbp</i>	Points to a <b>cfgncb</b> configuration notification control block.

On a platform that supports storage keys, the passed in *cbp* parameter must only be in the **KKEY\_PUBLIC** domain.

## Description

The **cfgnadd** kernel service adds a **cfgncb** control block to the list of **cfgncb** structures that the kernel maintains. A **cfgncb** control block contains the address of a notification routine (in its **cfgncb.func** field) to be called when a configurable variable is being changed.

The **SYS\_SETPARMS sysconfig** operation allows a user with sufficient authority to change the values of configurable system parameters. The **cfgnadd** service allows kernel routines and extensions to register the notification routine that is called whenever these configurable system variables have been changed.

This notification routine is called in a two-pass process. The first pass performs validity checks on the proposed changes to the system parameters. During the second pass invocation, the notification routine

performs whatever processing is needed to make these changes to the parameters. This two-pass procedure ensures that variables used by more than one kernel extension are correctly handled.

To use the **cfgnadd** service, the caller must define a **cfgncb** control block using the structure found in the **/usr/include/sys/sysconfig.h** file.

## Execution Environment

The **cfgnadd** kernel service can be called from the process environment only.

The **cfgncb.func** notification routine is called in a process environment only.

### Related reference:

“**cfgndel** Kernel Service” on page 40

### Related information:

**sysconfig** subroutine

Kernel Extension and Device Driver Management Kernel Services

## cfgncb Configuration Notification Control Block Purpose

Contains the address of a notification routine that is invoked each time the **sysconfig** subroutine is called with the **SYS\_SETPARMS** command.

## Syntax

```
int func (cmd, cur, new)
int cmd;
struct var *cur;
struct var *new;
```

## Parameters

Item	Description
<i>cmd</i>	Indicates the current operation type. Possible values are <b>CFGV_PREPARE</b> and <b>CFGV_COMMIT</b> , as defined in the <b>/usr/include/sys/sysconfig.h</b> file.
<i>cur</i>	Points to a <b>var</b> structure representing the current values of system-configurable variables.
<i>new</i>	Points to a <b>var</b> structure representing the new or proposed values of system-configurable variables.

The *cur* and *new* **var** structures are both in the system address space.

## Description

The configuration notification control block contains the address of a notification routine. This structure is intended to be used as a list element in a list of similar control blocks maintained by the kernel.

Each control block has the following definition:

```
struct cfgncb {
    struct cfgncb *cbnext;    /* next block on chain */
    struct cfgncb *cbprev;    /* prev control block on chain */
    int (*func)();           /* notification function */
};
```

The **cfgndel** or **cfgnadd** kernel service can be used to add or delete a **cfgncb** control block from the **cfgncb** list. To use either of these kernel services, the calling routine must define the **cfgncb** control block. This definition can be done using the **/usr/include/sys/sysconfig.h** file.

Every time a **SYS\_SETPARMS** **sysconfig** command is issued, the **sysconfig** subroutine iterates through the kernel list of **cfgncb** blocks, invoking each notification routine with a **CFGV\_PREPARE** command. This call represents the first pass of what is for the notification routine a two-pass process.

On a **CFGV\_PREPARE** command, the **cfgncb.func** notification routine should determine if any values of interest have changed. All changed values should be checked for validity. If the values are valid, a return code of 0 should be returned. Otherwise, a return value indicating the byte offset of the first field in error in the *new var* structure should be returned.

If all registered notification routines create a return code of 0, then no value errors have been detected during validity checking. In this case, the **sysconfig** subroutine issues its second pass call to the **cfgncb.func** routine and sends the same parameters, although the *cmd* parameter contains a value of **CFGV\_COMMIT**. This indicates that the new values go into effect at the earliest opportunity.

An example of notification routine processing might be the following. Suppose the user wishes to increase the size of the block I/O buffer cache. On a **CFGV\_PREPARE** command, the block I/O notification routine would verify that the proposed new size for the cache is legal. On a **CFGV\_COMMIT** command, the notification routine would then make the additional buffers available to the user by chaining more buffers onto the existing list of buffers.

**Related reference:**

“*cfgndel* Kernel Service”

**Related information:**

SYS\_SETPARMS subroutine

Kernel Extension and Device Driver Management Kernel Services

## **cfgndel Kernel Service**

### **Purpose**

Removes a notification routine for receiving broadcasts of changes to configurable system variables.

### **Syntax**

```
#include <sys/types.h> #include <sys/errno.h> #include <sys/sysconfig.h> void cfgndel ( cbp) struct  
cfgncb *cbp;
```

### **Parameter**

<b>Item</b>	<b>Description</b>
-------------	--------------------

<i>cbp</i>	Points to a <b>cfgncb</b> configuration notification control block.
------------	---

On a platform that supports storage keys, the passed in *cbp* parameter must only be in the **KKEY\_PUBLIC** domain.

### **Description**

The **cfgndel** kernel service removes a previously registered **cfgncb** configuration notification control block from the list of **cfgncb** structures maintained by the kernel. This service thus allows kernel routines and extensions to remove their notification routines from the list of those called when a configurable system variable has been changed.

The address of the **cfgncb** structure passed to the **cfgndel** kernel service must be the same address used to call the **cfgnadd** service when the structure was originally added to the list. The **/usr/include/sys/sysconfig.h** file contains a definition of the **cfgncb** structure.



## Execution Environment

The `cfgndel` kernel service can be called from the process environment only.

## Return Values

The `cfgndel` service has no return values.

### Related reference:

“`cfgnadd` Kernel Service” on page 38

### Related information:

`sysconfig` subroutine

Kernel Extension and Device Driver Management Kernel Services

## check Device Queue Management Routine

### Purpose

Provides a means for performing device-specific validity checking for parameters included in request queue elements.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/deviceq.h>
```

```
int check ( type, ptr, length)
int type;
struct req_qe *ptr;
int length;
```

### Parameters

Item	Description
<i>type</i>	Specifies the type of call. The following values are used when the kernel calls the <code>check</code> routine:  <code>CHECK_PARMS + SEND_CMD</code> Send command queue element.  <code>CHECK_PARMS + START_IO</code> Start I/O CCB queue element.  <code>CHECK_PARMS + GEN_PURPOSE</code> General purpose queue element.
<i>ptr</i>	Specifies the address of the queue element.
<i>length</i>	Specifies the length of the queue element.

### Description

The `check` routine is part of the Device Queue Management Kernel extension. Each device queue can have a `check` routine. This routine is optional and must be specified when the device queue is created with the `creatq` service. The `enqueue` service calls the `check` routine before a request queue element is put on the device queue. The kernel uses the routine's return value to determine whether to put the queue element on the device queue or to stop the request.

The kernel does not call the `check` routine when an acknowledgment or control queue element is sent. Therefore, the `check` routine is only called while executing within a process.

The address of the actual queue element is passed to this routine. In the **check** routine, take care to alter only the fields that were meant to be altered. This routine does not need to be serialized with the rest of the server's routines, because it is only checking the parameters in the queue element.

The **check** routine can check the request before the request queue element is placed on the device queue. The advantage of using this routine is that you can filter out unacceptable commands before they are put on the device queue.

The routine looks at the queue element and returns **RC\_GOOD** if the request is acceptable. If the return code is not **RC\_GOOD**, the kernel does not place the queue element in a device queue.

## Execution Environment

The **check** routine executes under the process environment of the requester. Therefore, access to data areas must be handled as if the routine were in an interrupt handler environment. There is, however, no requirement to pin the code and data as in a normal interrupt handler environment.

## Return Values

Item	Description
<b>RC_GOOD</b>	Indicates successful completion.

All other return values are device-specific.

### Related reference:

“enqueue Kernel Service” on page 137

## clrbuf Kernel Service

### Purpose

Sets the memory for the specified buffer structure's buffer to all zeros.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
void clrbuf ( bp)
struct buf *bp;
```

### Parameter

Item	Description
<i>bp</i>	Specifies the address of the buffer structure for the buffer to be cleared.

On a platform that supports storage keys, the passed in *bp* parameter must be in the **KKEY\_PUBLIC** or **KKEY\_BLOCK\_DEV** protection domain.

### Description

The **clrbuf** kernel service clears the buffer associated with the specified buffer structure. The **clrbuf** service does this by setting to 0 the memory for the buffer that contains the specified buffer structure.

## Execution Environment

The **clrbuf** kernel service can be called from either the process or interrupt environment.

## Return Values

The `clrbuf` service has no return values.

### Related information:

Block I/O Buffer Cache Kernel Services: Overview  
I/O Kernel Services

## clrjmpx Kernel Service Purpose

Removes a saved context by popping the last saved jump buffer from the list of saved contexts.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
void clrjmpx ( jump_buffer)
label_t *jump_buffer;
```

### Parameter

Item	Description
<i>jump_buffer</i>	Specifies the address of the caller-supplied jump buffer that was specified on the call to the <code>setjmpx</code> service.

### Description

The `clrjmpx` kernel service pops the most recent context saved by a call to the `setjmpx` kernel service. Since each `longjmpx` call automatically pops the jump buffer for the context to resume, the `clrjmpx` kernel service should be called only following:

- A normal return from the `setjmpx` service when the saved context is no longer needed
- Any code to be run that requires the saved context to be correct

The `clrjmpx` service takes the address of the jump buffer passed in the corresponding `setjmpx` service.

### Execution Environment

The `clrjmpx` kernel service can be called from either the process or interrupt environment.

### Return Values

The `clrjmpx` service has no return values.

### Related reference:

“setjmpx Kernel Service” on page 468

### Related information:

Process and Exception Management Kernel Services  
Understanding Exception Handling

## common\_reclock Kernel Service Purpose

Implements a generic interface to the record locking functions.

## Syntax

```
#include <sys/types.h>
#include <sys/flock.h>

common_relock( gp, size, offset,
lckdat, cmd, retray_fcn, retry_id, lock_fcn,
rele_fcn)
struct gnode *gp;
offset_t size;
offset_t offset;
struct eflock *lckdat;
int cmd;
int (*retray_fcn)();
ulong *retry_id;
int (*lock_fcn)();
int (*rele_fcn)();
```

## Parameters

Item	Description
<i>gp</i>	Points to the gnode that represents the file to lock.
<i>size</i>	Identifies the current size of the file in bytes.
<i>offset</i>	Specifies the current file offset. The system uses the <i>offset</i> parameter to establish where the lock region is to begin.
<i>lckdat</i>	Points to an <b>eflock</b> structure that describes the lock operation to perform.
<i>cmd</i>	Defines the type of operation the kernel service performs. This parameter is a bit mask consisting of the following bits:  <b>SETFLCK</b> If set, the system sets or clears a lock. If not set, the lock information is returned.  <b>SLPFLCK</b> If the lock cannot be granted immediately, wait for it. This is only valid when <b>SETFLCK</b> flag is set.  <b>INOFLCK</b> The caller is holding a lock on the object referred to by the gnode. The <b>common_relock</b> kernel service calls the release function before sleeping, and the lock function on return from sleep.  When the <i>cmd</i> parameter is set to <b>SLPFLCK</b> , it indicates that if the lock cannot be granted immediately, the service should wait for it. If the <i>retray_fcn</i> parameter contains a valid pointer, the <b>common_relock</b> kernel service does not sleep, regardless of the <b>SLPFLCK</b> flag.
<i>retray_fcn</i>	Points to a retry function. This function is called when the lock is retried. The retry function is not used if the lock is granted immediately. When the requested lock is blocked by an existing lock, a sleeping lock is established with the retry function address stored in it. The <b>common_relock</b> kernel service then returns a correlating ID (see the <i>retry_id</i> parameter) to the calling routine, along with an exit value of <b>EAGAIN</b> . When the sleeping lock is awakened, the retry function is called with the correlating ID as its ID argument.  If this argument is not <b>NULL</b> , then the <b>common_relock</b> kernel service does not sleep, regardless of the <b>SLPFLCK</b> command flag.
<i>retry_id</i>	Points to location to store the correlating ID. This ID is used to correlate a retry operation with a specific lock or set of locks. This parameter is used only in conjunction with retry function. The value stored in this location is an opaque value. The caller should not use this value for any purpose other than lock correlation.
<i>lock_fcn</i>	Points to a lock function. This function is invoked by the <b>common_relock</b> kernel service to lock a data structure used by the caller. Typically this is the data structure containing the gnode to lock. This function is necessary to serialize access to the object to lock. When the <b>common_relock</b> kernel service invokes the lock function, it is passed the private data pointer from the gnode as its only argument.
<i>rele_fcn</i>	Points to a release function. This function releases the lock acquired with the lock function. When the release function is invoked, it is passed the private data pointer from the gnode as its only argument.

## Description

The **common\_relock** routine implements a generic interface to the record-locking functions. This service allows distributed file systems to use byte-range locking. The kernel service does the following when a requested lock is blocked by an existing lock:

- Establishes a sleeping lock with the retry function in the **lock** structure. The address of the retry function is specified by the *retry\_fcn* parameter.
- Returns a correlating ID value to the caller along with an exit value of **EAGAIN**. The ID is stored in the *retry\_id* parameter.
- Calls the retry function when the sleeping lock is later awakened, the retry function is called with the *retry\_id* parameter as its argument.

**Note:** Before a call to the **common\_relock** subroutine, the **eflock** structure must be completely filled in. The *lckdat* parameter points to the **eflock** structure.

The caller can hold a serialization lock on the data object pointed to by the *gnode*. However, if the caller expects to sleep for a blocking-file lock and is holding the object lock, the caller must specify a lock function with the *lock\_fcn* parameter and a release function with the *rele\_fcn* parameter.

The lock is described by a **eflock** structure. This structure is identified by the *lckdat* parameter. If a read lock (**F\_RDLOCK**) or write lock (**F\_WRLOCK**) is set with a length of 0, the entire file is locked. Similarly, if unlock (**F\_UNLOCK**) is set starting at 0 for 0 length, all locks on this file are unlocked. This method is how locks are removed when a file is closed.

To allow the **common\_relock** kernel service to update the per-gnode lock list, the service takes a **GN\_RECLK\_LOCK** lock during processing.

## Execution Environment

The **common\_relock** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates successful completion.
<b>EAGAIN</b>	Indicates a lock cannot be granted because of a blocking lock and the caller did not request that the operation sleep.
<b>ERRNO</b>	Indicates an error. Refer to the <b>fcntl</b> system call for the list of possible values.

### Related information:

*fcntl* subroutine

*flock.h* subroutine

## compare\_and\_swap Kernel Services

### Purpose

Conditionally updates or returns a variable atomically.

### Syntax

```
#include <sys/atomic_op.h>
```

```
boolean_t compare_and_swap ( addr, old_val_addr, new_val)
atomic_p addr;
int * old_val_addr;
int new_val;
```

```
boolean_t compare_and_swaplp ( addr, old_val_addr, new_val)
atomic_l addr;
long * old_val_addr;
long new_val;
```

## Parameters

Item	Description
<i>addr</i>	Specifies the address of the variable.
<i>old_val_addr</i>	Specifies the address of the old value to be checked against (and conditionally updated with) the value of the variable.
<i>new_val</i>	Specifies the new value to be conditionally assigned to the variable.

## Description

The **compare\_and\_swap** kernel services performs an atomic (uninterruptible) operation which compares the contents of a variable with a stored old value; if equal, a new value is stored in the variable, and **TRUE** is returned, otherwise the old value is set to the current value of the variable, and **FALSE** is returned.

The **compare\_and\_swap** kernel service operates on a single word (32 bit) variable while the **compare\_and\_swaplp** kernel service operates on a double word (64 bit) variable.

The **compare\_and\_swap** kernel services are particularly useful in operations on singly linked lists, where a list pointer must not be updated if it has been changed by another thread since it was read.

### Note:

- The single word variable passed to the **compare\_and\_swap** kernel service must be aligned on a full word (32 bit) boundary.
- The double word variable passed to the **compare\_and\_swaplp** kernel service must be aligned on a double word (64 bit) boundary.

## Execution Environment

The **compare\_and\_swap** kernel services can be called from either the process or interrupt environment.

## Return Values

Item	Description
<b>TRUE</b>	Indicates that the variable was equal to the old value, and has been set to the new value.
<b>FALSE</b>	Indicates that the variable was not equal to the old value, and that its current value has been returned in the location where the old value was stored.

### Related reference:

“fetch\_and\_add Kernel Services” on page 140

“fetch\_and\_and or fetch\_and\_or Kernel Services” on page 141

### Related information:

Locking Kernel Services

## **coprocessor\_user\_register** Kernel Service Purpose

Registers the current process as a coprocessor user.

## Syntax

```
#include <sys/coprocessor.h>
kernno_t coprocessor_user_register ( int coprocessor_type, unsigned int * phandle )
```

## Parameters

Item	Description
<i>coprocessor_type</i>	Numeric value in the [0..63] range
<i>phandle</i>	Pointer to an unsigned 32 bit integer where a handle identifying this process is returned.

## Description

This kernel service allows a kernel extension to register the current process as a user of the coprocessor type passed as the first argument. When successful, the service sets up values in the process context that allow the current process to access coprocessors of the specified type in user mode.

## Execution Environment

This kernel service can be called in the process environment only.

## Return Values

When the call is successful, the kernel service returns a value of zero. Otherwise, a negative value is returned to indicate an error.

## Error Values

Possible errors are:

- Coprocessors not supported (supported only on POWER7<sup>®</sup> and newer processors)
- Invalid coprocessor type (must be in the range 0-63).
- Bad address passed as the second argument.
- The current process is already registered for this coprocessor type.
- The service is being called in interrupt context.
- The service could not allocate a value for the handle.

### Related information:

`coprocessor_user_unregister` subroutine

## `coprocessor_user_unregister` Kernel Service Purpose

Unregisters the current process as a coprocessor user.

## Syntax

```
#include <sys/coprocessor.h>
kernno_t coprocessor_user_unregister ( int coprocessor_type )
```

## Parameters

Item	Description
<i>coprocessor_type</i>	Numeric value in the range [0..63] which identifies a coprocessor type.

## Description

This kernel service allows a kernel extension to unregister the current process that was previously registered as a coprocessor user. When successful, further accesses by the process to the coprocessor type passed as an argument in user mode will fail with a privileged operation exception.

## Execution Environment

This kernel service can be called in the process environment only.

## Return Values

When the call is successful, the kernel service returns a value of zero. Otherwise, a negative value is returned to indicate an error.

## Error Values

Possible errors are:

- Coprocessors not supported (supported only on POWER7 and newer processors).
- Invalid coprocessor type (must be in the range 0-63).
- The current process is not registered for this coprocessor type.
- The service is being called in interrupt context.

### Related information:

`coprocessor_user_register` subroutine

## copyin Kernel Service

### Purpose

Copies data between user and kernel memory.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
int copyin ( uaddr, kaddr, count)
char *uaddr;
char *kaddr;
int count;
```

### Parameters



Item	Description
<i>uaddr</i>	Specifies the address of user data.
<i>kaddr</i>	Specifies the address of kernel data.
<i>count</i>	Specifies the number of bytes to copy.

## Description

The **copyin** kernel service copies the specified number of bytes from user memory to kernel memory. This service is provided so that system calls and device driver top half routines can safely access user data. The **copyin** service ensures that the user has the appropriate authority to access the data. It also provides recovery from paging I/O errors that would otherwise cause the system to crash.

The **copyin** service should be called only while executing in kernel mode in the user process.

## Execution Environment

The **copyin** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
EFAULT	Indicates that the user has insufficient authority to access the data, or the address specified in the <i>uaddr</i> parameter is not valid.
EIO	Indicates that a permanent I/O error occurred while referencing data.
ENOMEM	Indicates insufficient memory for the required paging operation.
ENOSPC	Indicates insufficient file system or paging space.

### Related reference:

“copyinstr Kernel Service”

“copyout Kernel Service” on page 50

### Related information:

Accessing User-Mode Data While in Kernel Mode

## copyinstr Kernel Service

### Purpose

Copies a character string (including the terminating null character) from user to kernel space.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/uio.h>
```

On the 32-bit kernel, the syntax for the **copyinstr** Kernel Service is:

```
int copyinstr (from, to, max, actual)
caddr_t from;
caddr_t to;
uint max;
uint *actual;
```

On the 64-bit kernel, the syntax for the **copyinstr** subroutine is:

```
int copyinstr (from, to, max, actual)
void *from;
void *to;
size_t max;
size_t *actual;
```

## Parameters

Item	Description
<i>from</i>	Specifies the address of the character string to copy.
<i>to</i>	Specifies the address to which the character string is to be copied.
<i>max</i>	Specifies the number of characters to be copied.
<i>actual</i>	Specifies a parameter, passed by reference, that is updated by the <b>copyinstr</b> service with the actual number of characters copied.

## Description

The **copyinstr** kernel service permits a user to copy character data from one location to another. The source location must be in user space or can be in kernel space if the caller is a kernel process. The destination is in kernel space.

## Execution Environment

The **copyinstr** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
E2BIG	Indicates insufficient space to complete the copy.
EIO	Indicates that a permanent I/O error occurred while referencing data.
ENOSPC	Indicates insufficient file system or paging space.
EFAULT	Indicates that the user has insufficient authority to access the data or the address specified in the <i>uaddr</i> parameter is not valid.

### Related information:

Accessing User-Mode Data While in Kernel Mode

Memory Kernel Services

## copyout Kernel Service

### Purpose

Copies data between user and kernel memory.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
int copyout ( kaddr, uaddr, count)
char *kaddr;
char *uaddr;
int count;
```

## Parameters

Item	Description
<i>kaddr</i>	Specifies the address of kernel data.
<i>uaddr</i>	Specifies the address of user data.
<i>count</i>	Specifies the number of bytes to copy.

## Description

The **copyout** service copies the specified number of bytes from kernel memory to user memory. It is provided so that system calls and device driver top half routines can safely access user data. The **copyout** service ensures that the user has the appropriate authority to access the data. This service also provides recovery from paging I/O errors that would otherwise cause the system to crash.

The **copyout** service should be called only while executing in kernel mode in the user process.

## Execution Environment

The **copyout** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
EFAULT	Indicates that the user has insufficient authority to access the data or the address specified in the <i>uaddr</i> parameter is not valid.
EIO	Indicates that a permanent I/O error occurred while referencing data.
ENOMEM	Indicates insufficient memory for the required paging operation.
ENOSPC	Indicates insufficient file system or paging space.

### Related reference:

“copyin Kernel Service” on page 48

“copyinstr Kernel Service” on page 49

### Related information:

Memory Kernel Services

## cpu\_speculation\_barrier kernel service

### Purpose

Provides protection against speculative execution side-channel attacks.

### Syntax

```
#include <sys/processor.h>
```

```
void cpu_speculation_barrier ( void )
```

### Description

The `cpu_speculation_barrier` kernel service provides kernel extensions with processor-model-dependent mitigation against known speculative-execution vulnerabilities. The `cpu_speculation_barrier` kernel service can be used to protect against side-channel attacks within the kernel environment. Kernel extensions should be carefully vetted when the `cpu_speculation_barrier` kernel service is used.

**Note:** Kernel performance might reduce when the `cpu_speculation_barrier` kernel service is used.

The `cpu_speculation_barrier` kernel service must be called before storage is accessed by using addresses that are computed from an untrusted source. Therefore, only kernel extensions that reference user-mode data directly without using cross-privilege domain access services, such as the `copyin` service, can use the `cpu_speculation_barrier` kernel service.

## Execution Environment

The `cpu_speculation_barrier` kernel service can be called from either the process environment or the interrupt environment.

## Example

The following example shows an `ioctl` device driver handler that directly references user-mode data:

```
int
dd_ioctl(dev_t devno, int cmd, void *arg, ulong devflag, chan_t chan, int ext)
{
    int    index;
    char   val;
    vector_t *uvec = NULL;
    extern int max_kdata_index;
    extern char kdata[];

    if (cmd == 0xC1C2) {
        /* Select kernel data from user input */
        uvec = (vector_t *)arg;
        index = uvec->index;

        if (index < max_kdata_index) {
            cpu_speculation_barrier();
            val = kdata[index];
            uvec->data[val]++;
        }
    }
}
```

## Return Values

The `cpu_speculation_barrier` kernel service does not return any value.

### Related information:

Accessing User-Mode Data While in Kernel Mode

## crCOPY Kernel Service

### Purpose

Copies a credentials structure to a new one and frees the old one.

### Syntax

```
#include <sys/cred.h>
```

```
struct ucred * crCOPY ( cr)
struct ucred * cr;
```

### Parameter

Item	Description
<i>cr</i>	Pointer to the credentials structure that is to be copied and then freed.

## Description

The **crCOPY** kernel service allocates a new credentials structure that is initialized from the contents of the *cr* parameter. The reference to *cr* is then freed and a pointer to the new structure returned to the caller.

**Note:** The *cr* parameter must have been obtained by an earlier call to the **crCOPY** kernel service, **crDUP** kernel service, **crGET** kernel service, or the **crREF** kernel service.

## Execution Environment

The **crCOPY** kernel service can be called from the process environment only.

## Return Values

Item	Description
Nonzero value	A pointer to a newly allocated and initialized credentials structure.
Zero value	An error occurred when the kernel service was attempting to allocate pinned memory for the credentials structure.

### Related information:

Security Kernel Services

## crDUP Kernel Service Purpose

Copies a credentials structure to a new one.

## Syntax

```
#include <sys/cred.h>
```

```
struct ucred * crdup ( cr)
struct ucred * cr;
```

## Parameter

Item	Description
<i>cr</i>	Pointer to the credentials structure that is to be copied.

## Description

The **crDUP** kernel service allocates a new credentials structure that is initialized from the contents of the *cr* parameter.

## Execution Environment

The **crDUP** kernel service can be called from the process environment only.

## Return Values

Item	Description
Nonzero value	A pointer to a newly allocated and initialized credentials structure.
Zero value	An error occurred when the kernel service was attempting to allocate pinned memory for the credentials structure.

### Related information:

Security Kernel Services

## creatp Kernel Service Purpose

Creates a new kernel process.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
pid_t creatp()
```

### Description

The **creatp** kernel service creates a kernel process. It also allocates and initializes a process block for the new process. Initialization involves these three tasks:

- Assigning an identifier to the kernel process.
- Setting the process state to idle.
- Initializing its parent, child, and sibling relationships.

"Using Kernel Processes" in *Kernel Extensions and Device Support Programming Concepts* has a more detailed discussion of how the **creatp** kernel service creates and initializes kernel processes.

The process calling the **creatp** service must subsequently call the **initp** kernel service to complete the process initialization. The **initp** service also makes the newly created process runnable.

### Execution Environment

The **creatp** kernel service can be called from the process environment only.

## Return Values

Item	Description
-1	Indicates an error.

Upon successful completion, the **creatp** kernel service returns the process identifier for the new kernel process.

### Related reference:

"initp Kernel Service" on page 221

### Related information:

Process and Exception Management Kernel Services

Using Kernel Processes

## CRED\_GETEUID, CRED\_GETRUID, CRED\_GETSUID, CRED\_GETLUID, CRED\_GETEGID, CRED\_GETRGID, CRED\_GETSGID and CRED\_GETNGRPS Macros

### Purpose

Credentials structure field accessing macros.

### Syntax

```
#include <sys/cred.h>
```

```
uid_t CRED_GETEUID ( crp )  
uid_t CRED_GETRUID ( crp )  
uid_t CRED_GETSUID ( crp )  
uid_t CRED_GETLUID ( crp )  
gid_t CRED_GETEGID ( crp )  
gid_t CRED_GETRGID ( crp )  
gid_t CRED_GETSGID ( crp )  
int CRED_GETNGRPS ( crp )
```

### Parameter

Item	Description
<i>crp</i>	Pointer to a credentials structure

### Description

These macros provide a means for accessing the user and group identifier fields within a credentials structure. The fields within a **ucred** structure should not be accessed directly as the field names and their locations are subject to change.

The **CRED\_GETEUID** macro returns the effective user ID field from the credentials structure referenced by *crp*.

The **CRED\_GETRUID** macro returns the real user ID field from the credentials structure referenced by *crp*.

The **CRED\_GETSUID** macro returns the saved user ID field from the credentials structure referenced by *crp*.

The **CRED\_GETLUID** macro returns the login user ID field from the credentials structure referenced by *crp*.

The **CRED\_GETEGID** macro returns the effective group ID field from the credentials structure referenced by *crp*.

The **CRED\_GETRGID** macro returns the real group ID field from the credentials structure referenced by *crp*.

The **CRED\_GETSGID** macro returns the saved group ID field from the credentials structure referenced by *crp*.

The **CRED\_GETNGRPS** macro returns the number of concurrent group ID values stored within the credentials structure referenced by *crp*.

These macros are defined in the system header file *<sys/cred.h>*.

## Execution Environment

The credentials macros called with any valid credentials pointer.

### Related information:

Security Kernel Services

## crexport Kernel Service

### Purpose

Copies an internal format credentials structure to an external format credentials structure.

### Syntax

```
#include <sys/cred.h>
```

```
void crexport (src, dst)
struct ucred * src;
struct ucred_ext * dst;
```

### Parameter

Item	Description
<i>src</i>	Pointer to the internal credentials structure.
<i>dst</i>	Pointer to the external credentials structure.

### Description

The **crexport** kernel service copies from the internal credentials structure referenced by *src* into the external credentials structure referenced by *dst*. The external credentials structure is guaranteed to be compatible between releases. Fields within a **ucred** structure must not be referenced directly as the field names and locations within that structure are subject to change.

## Execution Environment

The **crexport** kernel service can be called from the process environment only.

### Return Values

This kernel service does not have a return value.

### Related information:

Security Kernel Services

## crfree Kernel Service

### Purpose

Releases a reference count on a credentials structure.

### Syntax

```
#include <sys/cred.h>
```

```
void crfree ( cr)
struct ucred * cr;
```



## Parameter

Item	Description
<i>cr</i>	Pointer to the credentials structure that is to have a reference freed.

## Description

The **crfree** kernel service deallocates a reference to a credentials structure. The credentials structure is deallocated when no references remain.

**Note:** The *cr* parameter must have been obtained by an earlier call to the **crcopy** kernel service, **crdup** kernel service, **crget** kernel service, or the **crref** kernel service.

## Execution Environment

The **crfree** kernel service can be called from the process environment only.

## Return Values

No value is returned by this kernel service.

### Related information:

Security Kernel Services

## crget Kernel Service

### Purpose

Allocates a new, uninitialized credentials structure to a new one and frees the old one.

## Syntax

```
#include <sys/cred.h>
struct ucred * crget ( void )
```

## Parameter

This kernel service does not require any parameters.

## Description

The **crget** kernel service allocates a new credentials structure. The structure is initialized to all zero values, and the reference count is set to 1.

## Execution Environment

The **crget** kernel service can be called from the process environment only.

## Return Values

Item	Description
Nonzero value	A pointer to a newly allocated and initialized credentials structure.
Zero value	An error occurred when the kernel service was attempting to allocate pinned memory for the credentials structure.

**Related information:**

Security Kernel Services

**crhold Kernel Service**

**Purpose**

Increments the reference count for a credentials structure.

**Syntax**

```
#include <sys/cred.h>
```

```
void crhold ( cr)
struct ucred * cr;
```

**Parameter**

Item	Description
<i>cr</i>	Pointer to the credentials structure that will have its reference count incremented.

**Description**

The **crhold** kernel service increments the reference count of a credentials structure.

**Note:** Reference counts that are incremented with the **crhold** kernel service must be decremented with the **crfree** kernel service.

**Execution Environment**

The **crhold** kernel service can be called from the process environment only.

**Return Values**

No value is returned by this kernel service.

**Related information:**

Security Kernel Services

**crref Kernel Service**

**Purpose**

Increments the reference count for the current credentials structure.

**Syntax**

```
#include <sys/cred.h>
```

```
struct ucred * crref ( void )
```

## Parameter

This kernel service does not require any parameters.

## Description

The **crref** kernel service increments the reference count of the current credentials structure and returns a pointer to the current credentials structure to the invoker.

**Note:** References that are allocated with the **crref** kernel service must be released with the **crfree** kernel service.

## Execution Environment

The **crref** kernel service can be called from the process environment only.

## Return Values

Item	Description
Nonzero value	A pointer to the current credentials structure.

## Related information:

Security Kernel Services

## crset Kernel Service Purpose

Sets the current security credentials.

## Syntax

```
#include <sys/cred.h>
```

```
void crset ( cr )  
struct ucred * cr;
```

## Parameter

Item	Description
<i>cr</i>	Pointer to the credentials structure that will become the new, current security credentials.

## Description

The **crset** kernel service replaces the current security credentials with the supplied value. The existing structure will be deallocated.

**Note:** The *cr* parameter must have been obtained by an earlier call to the **crcopy** kernel service, **crdup** kernel service, **crget** kernel service, or the **crref** kernel service.

## Execution Environment

The **crset** kernel service can be called from the process environment only.

## Return Values

No value is returned by this kernel service.

### Related information:

Security Kernel Services

## curtime Kernel Service

### Purpose

Reads the current time into a time structure.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/time.h>
```

```
void curtime ( timestruct)
struct timestruc_t *timestruct;
```

### Parameter

Item	Description
<i>timestruct</i>	Points to a <b>timestruc_t</b> time structure defined in the <b>/usr/include/sys/time.h</b> file. The <b>curtime</b> kernel service updates the fields in this structure with the current time.

### Description

The **curtime** kernel service reads the current time into a time structure defined in the **/usr/include/sys/time.h** file. This service updates the **tv\_sec** and **tv\_nsec** fields in the time structure, pointed to by the *timestruct* parameter, from the hardware real-time clock. The kernel also maintains and updates a memory-mapped time **tod** structure. This structure is updated with each clock tick.

The kernel also maintains two other in-memory time values: the **lbolt** and **time** values. The three in-memory time values that the kernel maintains (the **tod**, **lbolt**, and **time** values) are available to kernel extensions. The **lbolt** in-memory time value is the number of timer ticks that have occurred since the system was booted. This value is updated once per timer tick. The **time** in-memory time value is the number of seconds since Epoch. The kernel updates the value once per second.

**Note:** POSIX 1003.1 defines "seconds since Epoch" as a "value interpreted as the number of seconds between a specified time and the Epoch". It further specifies that a "Coordinated Universal Time name specified in terms of seconds (*tm\_sec*), minutes (*tm\_min*), hours (*tm\_hour*), and days since January 1 of the year (*tm\_yday*), and calendar year minus 1900 (*tm\_year*) is related to a time represented as seconds since the Epoch, according to the following expression:  $tm\_sec + tm\_min * 60 + tm\_hour * 3600 + tm\_yday * 86400 + (tm\_year - 70) * 31536000 + ((tm\_year - 69) / 4) * 86400$  if the year is greater than or equal to 1970, otherwise it is undefined."

The **curtime** kernel service does not page-fault if a pinned stack and input time structure are used. Also, accessing the **lbolt**, **time**, and **tod** in-memory time values does not cause a page fault since they are in pinned memory.

### Execution Environment

The **curtime** kernel service can be called from either the process or interrupt environment.

The **tod**, **time**, and **lbolt** memory-mapped time values can also be read from the process or interrupt handler environment. The *timestruct* parameter and stack must be pinned when the **curtime** service is called in an interrupt handler environment.

## Return Values

The **curtime** kernel service has no return values.

### Related information:

Timer and Time-of-Day Kernel Services

## d

The following kernel services begin with the with the letter d.

### **d\_align** Kernel Service

#### Purpose

Provides needed information to align a buffer with a processor cache line.

#### Library

Kernel Extension Runtime Routines Library (**libsys.a**)

#### Syntax

```
int d_align()
```

#### Description

To maintain cache consistency with system memory, buffers must be aligned. The **d\_align** kernel service helps provide that function by returning the maximum processor cache-line size. The cache-line size is returned in log2 form.

#### Execution Environment

The **d\_align** service can be called from either the process or interrupt environment.

#### Related reference:

“**d\_cflush** Kernel Service” on page 62

“**d\_roundup** Kernel Service” on page 104

#### Related information:

Understanding Direct Memory Access (DMA) Transfer

### **d\_alloc\_dmamem** Kernel Service

#### Purpose

Allocates an area of “dma-able” memory.

#### Syntax

```
void * d_alloc_dmamem(d_handle_t device_handle, size_t size,int align)
```

#### Description

Exported, documented kernel service supported on PCI-based systems only. The **d\_alloc\_dmamem** kernel service allocates an area of “dma-able” memory which satisfies the constraints associated with a DMA handle, specified via the *device\_handle* parameter. The constraints (such as need for contiguous physical pages or need for 32-bit physical address) are intended to guarantee that a given adapter will be able to

access the physical pages associated with the allocated memory. A driver associates such constraints with a dma handle via the *flags* parameter on its **d\_map\_init** call.

The area to be allocated is the number of bytes in length specified by the *size* parameter, and is aligned on the byte boundary specified by the *align* parameter. The *align* parameter is actually the log base 2 of the desired address boundary. For example, an *align* value of 12 requests that the allocated area be aligned on a 4096 byte boundary.

**d\_alloc\_dmamem** is appropriate to be used for long-term mappings. Depending on the system configuration and the constraints encoded in the *device\_handle*, the underlying storage will come from either the *real\_heap* (**rmalloc** service) or *pinned\_heap* (**xmalloc** service).

**Note:**

1. The **d\_free\_dmamem** service should be called to free allocation from a previous **d\_alloc\_dmamem** call.
2. The **d\_alloc\_dmamem** kernel service can be called from the process environment only.

**Parameters**

Item	Description
<i>device_handle</i>	Indicates the dma handle.
<i>align</i>	Specifies alignment characteristics.
<i>size_t size</i>	Specifies number of bytes to allocate.

**Return Values**

Item	Description
Address of allocated area	Indicates that <b>d_alloc_dmamem</b> was successful.
NULL	Requested memory could not be allocated.

**Related reference:**

“d\_free\_dmamem Kernel Service” on page 76

“d\_map\_init Kernel Service” on page 82

“rmalloc Kernel Service” on page 449

**d\_cflush Kernel Service**

**Purpose**

Flushes the processor and I/O channel controller (IOCC) data caches when mapping bus device DMA with the long-term **DMA\_WRITE\_ONLY** option.

**Syntax**

```
int d_cflush (channel_id, baddr, count, daddr)
int channel_id;
caddr_t baddr;
size_t count;
caddr_t daddr;
```

**Parameters**

Item	Description
<i>channel_id</i>	Specifies the DMA channel ID returned by the <b>d_init</b> kernel service.
<i>baddr</i>	Designates the address of the memory buffer.
<i>count</i>	Specifies the length of the memory buffer transfer in bytes.
<i>daddr</i>	Designates the address of the device corresponding to the transfer.

## Description

The **d\_cflush** kernel service should be called after data has been modified in a buffer that will undergo direct memory access (DMA) processing. Through DMA processing, this data is sent to a device where the **d\_master** kernel service with the **DMA\_WRITE\_ONLY** option has already mapped the buffer for device DMA. The **d\_cflush** kernel service is not required if the **DMA\_WRITE\_ONLY** option is not used or if the buffer is mapped before each DMA operation by calling the **d\_master** kernel service.

The **d\_cflush** kernel service flushes the processor cache for the involved cache lines and invalidates any previously retrieved data that may be in the IOCC buffers for the designated channel. This most frequently occurs when using long-term buffer mapping for DMA support to or from a device.

### Long-Term DMA Buffer Mapping

The long-term DMA buffer mapping approach is frequently used when a pool of buffers is defined for sending commands and obtaining responses from an adapter using bus master DMA. This approach is also used frequently in the communications field where buffers can come from a common pool such as the **mbuf** pool or a pool used for protocol headers.

When using a fixed pool of buffers, the **d\_master** kernel service is used only once to map the pool's address and range. The device driver then modifies the data in the buffers. It must also flush the data from the processor and invalidate the IOCC data cache involved in transfers with the device. The IOCC cache must be invalidated because the data in the IOCC data cache may be stale due to the last DMA operation to or from the buffer area that has just been modified for the next operation.

The **d\_cflush** kernel service permits the flushing of the processor cache and making the required IOCC cache not valid. The device driver should use this service after modifying the data in the buffer and before sending the command to the device to start the DMA operation.

Once DMA processing has been completed, the device driver should call the **d\_complete** service to check for errors and ensure that any data read from the device has been flushed to memory.

**Note:** The **d\_cflush** kernel service is not supported on the 64-bit kernel.

## Execution Environment

The **d\_cflush** kernel service can be called from either the process or interrupt environment.

## Return Values

Item	Description
0	Indicates that the transfer was successfully completed.
EINVAL	Indicates the presence of an invalid parameter.

#### Related information:

I/O Kernel Services

Understanding Direct Memory Access (DMA) Transfer

## delay Kernel Service

### Purpose

Suspends the calling process for the specified number of timer ticks.

### Syntax

```
#include <sys/types.h>
```

```
#include <sys/errno.h>
```

```
void delay
```

```
( ticks)
```

```
int ticks;
```

### Parameter

Item	Description
<i>ticks</i>	Specifies the number of timer ticks that must occur before the process is reactivated. Many timer ticks can occur per second.

### Description

The **delay** kernel service suspends the calling process for the number of timer ticks specified by the *ticks* parameter.

The HZ value in the `/usr/include/sys/m_param.h` file can be used to determine the number of ticks per second.

### Execution Environment

The **delay** kernel service can be called from the process environment only.

### Return Values

The **delay** service has no return values.

#### Related information:

Timer and Time-of-Day Kernel Services

## del\_domain\_af Kernel Service

### Purpose

Deletes an address family from the Address Family domain switch table.

### Syntax

```
#include <sys/types.h>
```

```
#include <sys/errno.h>
```

```
#include <sys/domain.h>
```



```
int
del_domain_af ( domain)
struct domain *domain;
```

## Parameter

Item	Description
<i>domain</i>	Specifies the address family.

## Description

The `del_domain_af` kernel service deletes the address family specified by the *domain* parameter from the Address Family domain switch table.

## Execution Environment

The `del_domain_af` kernel service can be called from either the process or interrupt environment.

## Return Value

Item	Description
EINVAL	Indicates that the specified address is not found in the Address Family domain switch table.

## Example

To delete an address family from the Address Family domain switch table, invoke the `del_domain_af` kernel service as follows:

```
del_domain_af(&inetdomain);
```

In this example, the family to be deleted is `inetdomain`.

### Related reference:

“`add_domain_af` Kernel Service” on page 8

### Related information:

Network Kernel Services

## `del_input_type` Kernel Service

### Purpose

Deletes an input type from the Network Input table.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <net/if.h>
```

```
int del_input_type
( type)
u_short type;
```

### Parameter

Item	Description
<i>type</i>	Specifies which type of protocol the packet contains. This parameter is a field in a packet.

## Description

The **del\_input\_type** kernel service deletes an input type from the Network Input table to disable the reception of the specified packet type.

## Execution Environment

The **del\_input\_type** kernel service can be called from either the process or interrupt environment.

## Return Values

Item	Description
0	Indicates that the type was successfully deleted.
ENOENT	Indicates that the <b>del_input_type</b> service could not find the type in the Network Input table.

## Examples

1. To delete an input type from the Network Input table, invoke the **del\_input\_type** kernel service as follows:

```
del_input_type(ETHERTYPE_IP);
```

In this example, ETHERTYPE\_IP specifies that Ethernet IP packets should no longer be processed.

2. To delete an input type from the Network Input table, invoke the **del\_input\_type** kernel service as follows:

```
del_input_type(ETHERTYPE_ARP);
```

In this example, ETHERTYPE\_ARP specifies that Ethernet ARP packets should no longer be processed.

### Related reference:

“add\_input\_type Kernel Service” on page 8

“find\_input\_type Kernel Service” on page 143

### Related information:

Network Kernel Services

## del\_netisr Kernel Service

### Purpose

Deletes a network software interrupt service routine from the Network Interrupt table.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <net/netisr.h>
```

```
int del_netisr ( soft_intr_level)
u_short soft_intr_level;
```

### Parameter

Item	Description
<i>soft_intr_level</i>	Specifies the software interrupt level to delete. This parameter must be greater than or equal to 0 and less than <b>NETISR_MAX</b> . Refer to <b>netisr.h</b> for the range of values of <i>soft_intr_level</i> that are already in use. Also, other kernel extensions that are not AIX and that use network ISRs currently running on the system can make use of additional values not mentioned in <b>netisr.h</b> .

## Description

The **del\_netisr** kernel service deletes the network software interrupt service routine specified by the *soft\_intr\_level* parameter from the Network Software Interrupt table.

## Execution Environment

The **del\_netisr** kernel service can be called from either the process or interrupt environment.

## Return Values

Item	Description
0	Indicates that the software interrupt service was successfully deleted.
ENOENT	Indicates that the software interrupt service was not found in the Network Software Interrupt table.

## Example

To delete a software interrupt service from the Network Software Interrupt table, invoke the kernel service as follows:

```
del_netisr(NETISR_IP);
```

In this example, the software interrupt routine to be deleted is **NETISR\_IP**.

### Related reference:

“add\_netisr Kernel Service” on page 10

### Related information:

Network Kernel Services

## del\_netopt Macro

### Purpose

Deletes a network option structure from the list of network options.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <net/netopt.h>
```

```
del_netopt ( option_name_symbol )
option_name_symbol;
```

### Parameter

Item	Description
<i>option_name_symbol</i>	Specifies the symbol name used to construct the <b>netopt</b> structure and default names.

## Description

The **del\_netopt** macro deletes a network option from the linked list of network options. After the **del\_netopt** service is called, the option is no longer available to the **no** command.

## Execution Environment

The **del\_netopt** macro can be called from either the process or interrupt environment.

## Return Values

The **del\_netopt** macro has no return values.

### Related reference:

“add\_netopt Macro” on page 10

### Related information:

no command

Network Kernel Services

## detach Device Queue Management Routine

### Purpose

Provides a means for performing device-specific processing when the **detchq** kernel service is called.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/deviceq.h>
```

```
int detach( dev_parms, path_id)
caddr_t dev_parms;
cba_id path_id;
```

### Parameters

Item	Description
<i>dev_parms</i>	Passed to <b>creatd</b> service when the <b>detach</b> routine is defined.
<i>path_id</i>	Specifies the path identifier for the queue that is being detached from.

## Description

The **detach** routine is part of the Device Queue Management kernel extension. Each device queue can have a **detach** routine. This routine is optional and must be specified when the device queue is defined with the **creatd** service. The **detchq** service calls the **detach** routine each time a path to the device queue is removed.

To ensure that the **detach** routine is not called while a queue element from this client is still in the device queue, the kernel puts a detach control queue element at the end of the device queue. The server knows by convention that a detach control queue element signifies completion of all pending queue elements for that path. The kernel calls the **detach** routine after the detach control queue element is processed.

The **detach** routine executes under the process under which the **detachq** service is called. The kernel does not serialize the execution of this service with the execution of any of the other server routines.

## Execution Environment

The **detach** routine can be called from the process environment only.

## Return Values

Item	Description
RC_GOOD	Indicates successful completion.

A return value other than **RC\_GOOD** indicates an irrecoverable condition causing system failure.

## devdump Kernel Service Purpose

Calls a device driver dump-to-device routine.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
int devdump
(devno, uiop, cmd, arg, chan, ext)
dev_t devno;
struct uio * uiop;
int cmd, arg, ext;
```

## Parameters

Item	Description
<i>devno</i>	Specifies the major and minor device numbers.
<i>uiop</i>	Points to the <b>uio</b> structure containing write parameters.
<i>cmd</i>	Specifies which dump command to perform.
<i>arg</i>	Specifies a parameter or address to a parameter block for the specified command.
<i>chan</i>	Specifies the channel ID.
<i>ext</i>	Specifies the extended system call parameter.

## Description

The kernel or kernel extension calls the **devdump** kernel service to initiate a memory dump to a device when writing dump data and then to terminate the dump to the target device.

The **devdump** service calls the device driver's **dddump** routine, which is found in the device switch table for the device driver associated with the specified device number. If the device number (specified by the *devno* parameter) is not valid or if the associated device driver does not have a **dddump** routine, an **ENODEV** return value is returned.

If the device number is valid and the specified device driver has a **dddump** routine, the routine is called.

If the device driver's **dddump** routine is successfully called, the return value for the **devdump** service is set to the return value provided by the device's **dddump** routine.

## Execution Environment

The **devdump** kernel service can be called in either the process or interrupt environment, as described under the conditions described in the **dddump** routine.

## Return Values

Item	Description
0	Indicates a successful operation.
ENODEV	Indicates that the device number is not valid or that no <b>dddump</b> routine is registered for this device.

The **dddump** device driver routine provides other return values.

### Related reference:

“dddump Device Driver Entry Point” on page 624

### Related information:

Kernel Extension and Device Driver Management Kernel Services

## devstrat Kernel Service

### Purpose

Calls a block device driver's strategy routine.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
int devstrat ( bp)
struct buf *bp;
```

### Parameter

Item	Description
<i>bp</i>	Points to the <b>buf</b> structure specifying the block transfer parameters.

### Description

The kernel or kernel extension calls the **devstrat** kernel service to request a block data transfer to or from the device with the specified device number. This device number is found in the **buf** structure. The **devstrat** service can only be used for the block class of device drivers.

The **devstrat** service calls the device driver's **ddstrategy** routine. This routine is found in the device switch table for the device driver associated with the specified device number found in the **b\_dev** field. The **b\_dev** field is found in the **buf** structure pointed to by the *bp* parameter. The caller of the **devstrat** service must have an **iodone** routine specified in the **b\_iodone** field of the **buf** structure. Following the return from the device driver's **ddstrategy** routine, the **devstrat** service returns without waiting for the I/O to be performed.

On multiprocessor systems, all **iodone** routines run by default on the first processor started when the system was booted. This ensures compatibility with uniprocessor device drivers. If the **iodone** routine has been designed to be multiprocessor-safe, set the **B\_MPSAFE** flag in the **b\_flags** field of the **buf** structure passed to the **devstrat** kernel service. The **iodone** routine will then run on any available processor.

If the device major number is not valid or the specified device is not a block device driver, the **devstrat** service returns the **ENODEV** return code. If the device number is valid, the device driver's **ddstrategy**

routine is called with the pointer to the **buf** structure (specified by the *bp* parameter).

## Execution Environment

The **devstrat** kernel service can be called from either the process or interrupt environment.

**Note:** The **devstrat** kernel service can be called in the interrupt environment only if its priority level is INTIODONE or lower.

## Return Values

Item	Description
0	Indicates a successful operation.
ENODEV	Indicates that the device number is not valid or that no <b>ddstrategy</b> routine registered. This value is also returned when the specified device is not a block device driver. If this error occurs, the <b>devstrat</b> service can cause a page fault.

### Related reference:

“ddstrategy Device Driver Entry Point” on page 635

“buf Structure” on page 615

### Related information:

Kernel Extension and Device Driver Management Kernel Services

## devswadd Kernel Service

### Purpose

Adds a device entry to the device switch table.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/device.h>
```

```
int devswadd ( devno, dswptr)
dev_t devno;
struct devsw *dswptr;
```

### Parameters

Item	Description
<i>devno</i>	Specifies the major and minor device numbers to be associated with the specified entry in the device switch table.
<i>dswptr</i>	Points to the device switch structure to be added to the device switch table.

### Description

The **devswadd** kernel service is typically called by a device driver's **ddconfig** routine to add or replace the device driver's entry points in the device switch table. The device switch table is a table of device switch (**devsw**) structures indexed by the device driver's major device number. This table of structures is used by the device driver interface services in the kernel to facilitate calling device driver routines.

The major device number portion of the *devno* parameter is used to specify the index in the device switch table where the **devswadd** service must place the specified device switch entry. Before this service copies the device switch structure into the device switch table, it checks the existing entry to determine if any opened device is using it. If an opened device is currently occupying the entry to be replaced, the **devswadd** service does not perform the update. Instead, it returns an **EEXIST** error value. If the update is successful, it returns a value of 0.

Entry points in the device switch structure that are not supported by the device driver must be handled in one of two ways. If a call to an unsupported entry point should result in the return of an error code, then the entry point must be set to the **nodev** routine in the structure. As a result, any call to this entry point automatically invokes the **nodev** routine, which returns an **ENODEV** error code. The kernel provides the **nodev** routine.

Otherwise, a call to an unsupported entry point should be treated as a no-operation function. Then the corresponding entry point should be set to the **nulldev** routine. The **nulldev** routine, which is also provided by the kernel, performs no operation if called and returns a 0 return code.

On multiprocessor systems, all device driver routines run by default on the first processor started when the system was booted. This ensures compatibility with uniprocessor device drivers. If the device driver being added has been designed to be multiprocessor-safe, set the **DEV\_MPSAFE** flag in the **d\_opts** field of the **devsw** structure passed to the **devswadd** kernel service. The device driver routines will then run on any available processor.

All other fields within the structure that are not used should be set to 0. Some fields in the structure are for kernel use; the **devswadd** service does not copy these fields into the device switch table. These fields are documented in the **/usr/include/device.h** file.

## Execution Environment

The **devswadd** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
EEXIST	Indicates that the specified device switch entry is in use and cannot be replaced.
ENOMEM	Indicates that the entry cannot be pinned due to insufficient real memory.
EINVAL	Indicates that the major device number portion of the <i>devno</i> parameter exceeds the maximum permitted number of device switch entries.

### Related reference:

“devswdel Kernel Service” on page 73

“ddconfig Device Driver Entry Point” on page 621

### Related information:

Kernel Extension and Device Driver Management Kernel Services

## devswchg Kernel Service

### Purpose

Alters a device switch entry point in the device switch table.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/device.h>
```

```
int devswchg ( devno, type, newfunc, oldfunc);
dev_t devno;
int type;
int (*newfunc) ();
int (**oldfunc)();
```



## Parameters

Item	Description
<i>devno</i>	Specifies the major and minor device numbers of the device to be changed.
<i>type</i>	Specifies the device switch entry point to alter. The <i>type</i> parameter can have one of the following values: <b>DSW_BLOCK</b> Alters the <b>ddstrategy</b> entry point. <b>DSW_CONFIG</b> Alters the <b>ddconfig</b> entry point. <b>DSW_CREAD</b> Alters the <b>ddread</b> entry point. <b>DSW_CWRITE</b> Alters the <b>ddwrite</b> entry point. <b>DSW_DUMP</b> Alters the <b>dddump</b> entry point. <b>DSW_MPX</b> Alters the <b>ddmpx</b> entry point. <b>DSW_SELECT</b> Alters the <b>ddselect</b> entry point. <b>DSW_TCPATH</b> Alters the <b>ddrevoke</b> entry point.
<i>newfunc</i>	Specifies the new value for the device switch entry point.
<i>oldfunc</i>	Specifies that the old value of the device switch entry point be returned here.

## Description

The **devswchg** kernel service alters the value of a device switch entry point (function pointer) after a device switch table entry has been added by the **devswadd** kernel service. The device switch entry point specified by the *type* parameter is set to the value of the *newfunc* parameter. Its previous value is returned in the memory addressed by the *oldfunc* parameter. Only one device switch entry can be altered per call.

If the **devswchg** kernel service is unsuccessful, the value referenced by the *oldfunc* parameter is not defined.

## Execution Environment

The **devswchg** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
EINVAL	Indicates the <i>Type</i> command was not valid.
ENODEV	Indicates the device switch entry specified by the <i>devno</i> parameter is not defined.

### Related reference:

“devswadd Kernel Service” on page 71

### Related information:

List of Kernel Extension and Device Driver Management Kernel Services

## devswdel Kernel Service

### Purpose

Deletes a device driver entry from the device switch table.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/device.h>
```

```
int devswdel
( devno)
dev_t devno;
```

## Parameter

Item	Description
<i>devno</i>	Specifies the major and minor device numbers of the device to be deleted.

## Description

The **devswdel** kernel service is typically called by a device driver's **ddconfig** routine on termination to remove the device driver's entry points from the device switch table. The device switch table is a table of device switch (**devsw**) structures indexed by the device driver's major device number. The device driver interface services use this table of structures in the kernel to facilitate calling device driver routines.

The major device number portion of the *devno* parameter is used to specify the index into the device switch table for the entry to be removed. Before the device switch structure is removed, the existing entry is checked to determine if any opened device is using it.

If an opened device is currently occupying the entry to be removed, the **devswdel** service does not perform the update. Instead, it returns an **EEXIST** return code. If the removal is successful, a return code of 0 is set.

The **devswdel** service removes a device switch structure entry from the table by marking the entry as undefined and setting all of the entry point fields within the structure to a **nodev** value. As a result, any callers of the removed device driver return an **ENODEV** error code. If the specified entry is already marked undefined, the **devswdel** service returns an **ENODEV** error code.

## Execution Environment

The **devswdel** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
<b>EEXIST</b>	Indicates that the specified device switch entry is in use and cannot be removed.
<b>ENODEV</b>	Indicates that the specified device switch entry is not defined.
<b>EINVAL</b>	Indicates that the major device number portion of the <i>devno</i> parameter exceeds the maximum permitted number of device switch entries.

### Related reference:

“devswchg Kernel Service” on page 72

“devswqry Kernel Service” on page 75

### Related information:

Kernel Extension and Device Driver Management Kernel Services

## devswqry Kernel Service

### Purpose

Checks the status of a device switch entry in the device switch table.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/device.h>int devswqry ( devno, status, dsdptr)
dev_t devno;
uint *status;
caddr_t *dsdptr;
```

### Parameters

Item	Description
<i>devno</i>	Specifies the major and minor device numbers of the device to be queried.
<i>status</i>	Points to the status of the specified device entry in the device switch table. This parameter is passed by reference.
<i>dsdptr</i>	Points to device-dependent information for the specified device entry in the device switch table. This parameter is passed by reference.

### Description

The **devswqry** kernel service returns the status of a specified device entry in the device switch table. The entry in the table to query is determined by the major portion of the device number specified in the *devno* parameter. The status of the entry is returned in the *status* parameter that is passed by reference on the call. If this pointer is null on entry to the **devswqry** service, then the status is not returned to the caller.

The **devswqry** service also returns the address of device-dependent information for the specified device entry in the device switch table. This address is taken from the *d\_dsdptr* field for the entry and returned in the *dsdptr* parameter, which is passed by reference. If this pointer is null on entry to the **devswqry** service, then the service does not return the address from the *d\_dsdptr* field to the caller.

### Status Parameter Flags

The *status* parameter comprises a set of flags that can indicate the following conditions:

Item	Description
<b>DSW_BLOCK</b>	Device switch entry is defined by a block device driver. This flag is set when the device driver has a <b>ddstrategy</b> entry point.
<b>DSW_CONFIG</b>	Device driver in this device switch entry provides an entry point for configuration.
<b>DSW_CREAD</b>	Device driver in this device switch entry is providing a routine for character reads or raw input. This flag is set when the device driver has a <b>ddread</b> entry point.
<b>DSW_CWRITE</b>	Device driver in this device switch entry is providing a routine for character writes or raw output. This flag is set when the device driver has a <b>ddwrite</b> entry point.
<b>DSW_DEFINED</b>	Device switch entry is defined.
<b>DSW_DUMP</b>	Device driver defined by this device switch entry provides the capability to support one or more of its devices as targets for a kernel dump. This flag is set when the device driver has provided a <b>dddump</b> entry point.
<b>DSW_MPX</b>	Device switch entry is defined by a multiplexed device driver. This flag is set when the device driver has a <b>ddmpx</b> entry point.
<b>DSW_OPENED</b>	Device switch entry is in use and the device has outstanding opens. This flag is set when the device driver has at least one outstanding open.
<b>DSW_SELECT</b>	Device driver in this device switch entry provides a routine for handling the <b>select</b> or <b>poll</b> subroutines. This flag is set when the device driver has provided a <b>ddselect</b> entry point.

Item	Description
DSW_TCPATH	Device driver in this device switch entry supports devices that are considered to be in the trusted computing path and provide support for the revoke function. This flag is set when the device driver has provided a <b>ddrevoke</b> entry point.
DSW_TTY	Device switch entry is in use by a tty device driver. This flag is set when the pointer to the <b>d_ttys</b> structure is not a null character.
DSW_UNDEFINED	Device switch entry is not defined.

The *status* parameter is set to the **DSW\_UNDEFINED** flag when a device switch entry is not in use. This is the case if either of the following are true:

- The entry has never been used. (No previous call to the **devswadd** service was made.)
- The entry has been used but was later deleted. (A call to the **devswadd** service was issued, followed by a call to the **devswdel** service.)

No other flags are set when the **DSW\_UNDEFINED** flag is set.

**Note:** The *status* parameter must be a null character if called from the interrupt environment.

## Execution Environment

The **devswqry** kernel service can be called from either the process or interrupt environment.

## Return Values

Item	Description
0	Indicates a successful operation.
EINVAL	Indicates that the major device number portion of the <i>devno</i> parameter exceeds the maximum permitted number of device switch entries.

### Related reference:

“devswadd Kernel Service” on page 71

“devswchg Kernel Service” on page 72

### Related information:

Kernel Extension and Device Driver Management Kernel Services

## d\_free\_dmamem Kernel Service

### Purpose

Frees an area of memory.

### Syntax

```
int d_free_dmamem(d_handle_t device_handle, void * addr, size_t size)
```

### Description

Exported, documented kernel service supported on PCI-based systems only. The **d\_free\_dmamem** kernel service frees the area of memory pointed to by the *addr* parameter. This area of memory must be allocated with the **d\_alloc\_dmamem** kernel service using the same *device\_handle*, and the *addr* must be the address returned from the corresponding **d\_alloc\_dmamem** call. Also, the size must be the same size that was used on the corresponding **d\_alloc\_dmamem** call.

### Note:

1. Any memory allocated in a prior **d\_alloc\_dmamem** call must be explicitly freed with a **d\_free\_dmamem** call.

2. This service can be called from the process environment only.

## Parameters

Item	Description
<i>device_handle</i>	Indicates the dma handle.
<i>size_t size</i>	Specifies size of area to free.
<i>void * addr</i>	Specifies address of area to free.

## Return Values

Item	Description
0	Indicates successful completion.
-1	Indicates underlying free service (xmfree or rmalloc) failed.

### Related reference:

“d\_alloc\_dmamem Kernel Service” on page 61

## disable\_lock Kernel Service

### Purpose

Raises the interrupt priority, and locks a simple lock if necessary.

### Syntax

```
#include <sys/lock_def.h>

int disable_lock ( int_pri, lock_addr)
int int_pri;
simple_lock_t lock_addr;
```

## Parameters

Item	Description
<i>int_pri</i>	Specifies the interrupt priority to set.
<i>lock_addr</i>	Specifies the address of the lock word to lock.

## Description

The **disable\_lock** kernel service raises the interrupt priority, and locks a simple lock if necessary, in order to provide optimized thread-interrupt critical section protection for the system on which it is executing. On a multiprocessor system, calling the **disable\_lock** kernel service is equivalent to calling the **i\_disable** and **simple\_lock** kernel services. On a uniprocessor system, the call to the **simple\_lock** service is not necessary, and is omitted. However, you should still pass a valid lock address to the **disable\_lock** kernel service. Never pass a **NULL** lock address.

## Execution Environment

The **disable\_lock** kernel service can be called from either the process or interrupt environment.

## Return Values

The **disable\_lock** kernel service returns the previous interrupt priority.

### Related reference:

“i\_disable Kernel Service” on page 208

“simple\_lock\_init Kernel Service” on page 475

## Related information:

Understanding Locking

## disablement\_checking\_resume Kernel Service

### Purpose

Indicates the end of a disabled code path that was exempted from detection of excessive interrupt disablement.

### Syntax

```
#include <sys/intr.h>
```

```
void disablement_checking_resume(long prev_state)
```

### Parameters

Item	Description
<i>prev_state</i>	Specifies the disablement detection state to be restored. This value is returned by the <b>disablement_checking_suspend</b> kernel service.

### Description

The **disablement\_checking\_resume** service restores the disablement detection state to the value passed as *prev\_state*. This service must be called after reenabling interrupts at the end of an INTMAX critical section, not within it. This is because, in the case of an INTMAX critical section, the tick counting will have been deferred by the total disablement until the moment of enablement.

This service must be used in conjunction with the **disablement\_checking\_suspend** kernel service, which temporarily stops disablement detection.

**Note:** Error checking, including that for excessive interrupt disablement, can be enabled or disabled by the **errctrl** command.

### Execution Environment

The **disablement\_checking\_resume** service can be called from either the process or the interrupt environments.

#### Related reference:

“disablement\_checking\_suspend Kernel Service”

#### Related information:

errctrl command

## disablement\_checking\_suspend Kernel Service

### Purpose

Indicates the start of a disabled code path that is exempt from detection of excessive interrupt disablement.

### Syntax

```
#include <sys/intr.h>
```

```
long disablement_checking_suspend(void)
```

## Description

A call to the **disablement\_checking\_suspend** service temporarily disables the detection of excessive disablement for the duration of a portion of a critical section. For base level code, insert this call at the beginning of the exempt critical section immediately after it disables, or as soon as possible within interrupt handling code.

This service must be used in conjunction with the **disablement\_checking\_resume** kernel service, which resumes the prior disablement checking state.

**Note:** Error checking, including that for excessive interrupt disablement, can be enabled or disabled by the **errctrl** command.

## Execution Environment

The **disablement\_checking\_suspend** service can be called from either the process or the interrupt environments. Interrupts should be at least partially disabled at the time of the call.

## Return Values

The **disablement\_checking\_suspend** service returns the previous suspension state to the caller. This value must be passed later to the resume function, which restores that state. This enables nesting of exempt critical sections.

### Related reference:

“disablement\_checking\_resume Kernel Service” on page 78

### Related information:

errctrl command

## d\_map\_attr Kernel Service Purpose

Changes the attributes associated with a DMA handle.

## Syntax

```
#include <sys/dma.h>
kernno_t d_map_attr (handle, cmd, attr, attr_size)
d_handle_t handle;
ulong cmd;
void * attr;
size_t attr_size;
```

## Parameters

Item	Description
<i>handle</i>	Indicates the unique handle returned by the <b>d_map_init_ext</b> kernel service.
<i>cmd</i>	Specifies one of the following flags:  <b>D_ATTR_SET_MIN_MAPMEM</b> Sets the minimum amount of I/O mappable memory. This is the logical memory change and not the DMA bus memory change.  <b>D_ATTR_SET_DES_MAPMEM</b> Sets the desired amount of I/O mappable memory. This is the logical memory change and not the DMA bus memory change.
<i>attr</i>	You must set this parameter to the value of <b>size64_t *</b> . This parameter sets the minimum or the desired amount of I/O mappable memory depending on the specified value of the <i>cmd</i> parameter.
<i>attr_size</i>	You must set this parameter to the value of <b>sizeof(size64_t)</b> . This parameter sets the minimum or the desired amount of I/O mappable memory depending on the specified value of the <i>cmd</i> parameter.

## Description

The **d\_map\_attr** kernel service can change certain attributes of the **d\_handle\_t** structure in case the needs of a device driver change during runtime. For example, if a device driver needs more DMA space at runtime, it can call the **d\_map\_attr** kernel service to request an increase in the map space. The **d\_map\_attr** kernel service is not an exported kernel service, but a bus specific utility routine determined by the **d\_map\_init\_ext** kernel service and provided to the caller through the **d\_handle** structure.

## Execution Environment

The **d\_map\_attr** kernel service can be called from the process environment at **INTBASE**. Serialization with other DMA services like the **d\_map\_page** service and the **d\_unmap\_page** service is the caller's responsibility.

## Return Values

Item	Description
DMA_SUCC	Indicates a successful completion.
EINVAL_D_MAP_ATTR	Indicates that the specified <i>cmd</i> parameter is not valid.
ENOMEM_D_MAP_ATTR	Indicates that it is unable to change the minimum or desired I/O mappable memory.

### Related reference:

“d\_map\_init\_ext Kernel Service” on page 83

## d\_map\_clear Kernel Service

### Purpose

Deallocates resources previously allocated on a **d\_map\_init** call.

### Syntax

```
#include <sys/dma.h>
```

```
void d_map_clear (*handle)  
struct d_handle *handle
```

### Parameters

Item	Description
<i>handle</i>	Indicates the unique handle returned by the <b>d_map_init</b> kernel service.

## Description

The **d\_map\_clear** kernel service is a bus-specific utility routine determined by the **d\_map\_init** service that deallocates resources previously allocated on a **d\_map\_init** call. This includes freeing the **d\_handle** structure that was allocated by **d\_map\_init**.

**Note:** You can use the **D\_MAP\_CLEAR** macro provided in the **/usr/include/sys/dma.h** file to code calls to the **d\_map\_clear** kernel service.

### Related reference:

“d\_map\_init Kernel Service” on page 82

## d\_map\_disable Kernel Service

### Purpose

Disables DMA for the specified handle.



## Syntax

```
#include <sys/dma.h>
int d_map_disable(*handle)
struct d_handle *handle;
```

## Parameters

Item	Description
<i>handle</i>	Indicates the unique handle returned by <code>d_map_init</code> .

## Description

The `d_map_disable` kernel service is a bus-specific utility routine determined by the `d_map_init` kernel service that disables DMA for the specified *handle* with respect to the platform.

**Note:** You can use the `D_MAP_DISABLE` macro provided in the `/usr/include/sys/dma.h` file to code calls to the `d_map_disable` kernel service.

## Return Values

Item	Description
<code>DMA_SUCC</code>	Indicates the DMA is successfully disabled.
<code>DMA_FAIL</code>	Indicates the DMA could not be explicitly disabled for this device or bus.

## Related reference:

“`d_map_init` Kernel Service” on page 82

## `d_map_enable` Kernel Service

### Purpose

Enables DMA for the specified handle.

## Syntax

```
#include <sys/dma.h>
int d_map_enable(*handle)
struct d_handle *handle;
```

## Parameters

Item	Description
<i>handle</i>	Indicates the unique handle returned by <code>d_map_init</code> .

## Description

The `d_map_enable` kernel service is a bus-specific utility routine determined by the `d_map_init` kernel service that enables DMA for the specified *handle* with respect to the platform.

**Note:** You can use the `D_MAP_ENABLE` macro provided in the `/usr/include/sys/dma.h` file to code calls to the `d_map_enable` kernel service.

## Return Values

Item	Description
DMA_SUCC	Indicates the DMA is successfully enabled.
DMA_FAIL	Indicates the DMA could not be explicitly enabled for this device or bus.

#### Related reference:

“d\_map\_init Kernel Service”

## d\_map\_init Kernel Service

### Purpose

Allocates and initializes resources for performing DMA with PCI and ISA devices.

### Syntax

```
#include <sys/dma.h>
```

```
struct d_handle* d_map_init (bid, flags, bus_flags, channel)
int bid;
int flags;
int bus_flags;
uint channel;
```

### Parameters

Item	Description
<i>bid</i>	Specifies the bus identifier.
<i>flags</i>	Describes the mapping.
<i>bus_flags</i>	Specifies the target bus flags.
<i>channel</i>	Indicates the <i>channel</i> assignment specific to the bus.

### Description

The **d\_map\_init** kernel service allocates and initializes resources needed for managing DMA operations and returns a unique *handle* to be used on subsequent DMA service calls. The *handle* is a pointer to a **d\_handle** structure allocated by **d\_map\_init** from the pinned heap for the device. The device driver uses the function addresses provided in the *handle* for accessing the DMA services specific to its host bus. The **d\_map\_init** service returns a **DMA\_FAIL** error when resources are unavailable or cannot be allocated.

The *channel* parameter is the assigned channel number for the device, if any. Some devices and or buses might not have the concept of *channels*. For example, an ISA device driver would pass in its assigned DMA channel in the *channel* parameter.

**Note:** The possible flag values for the *flags* parameter can be found in `/usr/include/sys/dma.h`. These flags can be logically ORed together to reflect the desired characteristics.

### Execution Environment

The **d\_map\_init** kernel service should only be called from the process environment.

### Return Values

Item	Description
DMA_FAIL	Indicates that the resources are unavailable. No registration was completed.
struct d_handle *	Indicates successful completion.

#### Related reference:

“d\_unmap\_page Kernel Service” on page 108

“d\_map\_list Kernel Service” on page 84

“d\_map\_disable Kernel Service” on page 80

## d\_map\_init\_ext Kernel Service

### Purpose

Allocates and initializes resources for performing DMA with PCI and VDEVICE devices.

### Syntax

```
#include <sys/types.h>
#include <sys/dma.h>
#include <sys/kerrno.h>

kerrno_t d_map_init_ext (dma_input, info_size, handle_ptr)
d_info_t * dma_input;
size_t info_size;
d_handle_t * handle_ptr;
```

### Parameters

Item	Description
<i>dma_input</i>	Contains information like the bus identifier, flags, and so on.
<i>info_size</i>	Specifies the size of the <i>dma_input</i> parameter in bytes.
<i>handle_ptr</i>	Contains the DMA handle returned upon success.

### Description

The **d\_map\_init\_ext** kernel service is very similar to the **d\_map\_init** kernel service. Unlike the **d\_map\_init** kernel service, the input argument list of the **d\_map\_init\_ext** kernel service is not limited and can be extended without breaking binary compatibility. Also, the **d\_map\_init\_ext** kernel service returns a **kerrno\_t** type return code which contains more RAS information rather than just the **DMA\_FAIL** value.

The caller of the **d\_map\_init\_ext** kernel service initializes the **d\_info\_t** structure and passes it into the **d\_map\_init\_ext** kernel service by reference. The size of the **d\_info\_t** type must match the *info\_size* parameter. This allows future expansion of the **d\_info\_t** type safely. If there is a size mismatch, the **d\_map\_init\_ext** kernel service fails. The **d\_map\_init\_ext** kernel service also creates a new private pool of I/O memory entitlement that can be used for DMA. The private pool is created by carving out a chunk of total I/O memory entitlement for the AIX partition. Thus, in order to create a **d\_handle\_t** type successfully, there must be sufficient DMA PCI space and I/O memory entitlement.

The following structure is defined in the **sys/dma.h** file:

```
#define DMA_MAX_MAPPER_NAME    32
typedef struct d_info
{
    uint64_t    di_bid;
    uint64_t    di_flags;
    uint64_t    di_bus_flags;
    uint64_t    di_channel;
    uint64_t    di_min_mapmem;
```

```

uint64_t    di_des_mapmem;
uint64_t    di_max_mapmem;
char        di_mapper_name[DMA_MAX_MAPPER_NAME];
} d_info_t;

```

**Note:** The first four fields of the `d_info_t` type match the four arguments of the `d_map_init` kernel service. Therefore, all flags and bus\_flags on the `d_map_init` kernel service are honored by the `d_map_init_ext` kernel service except the `DMA_MAXMIN_*` flags. The `DMA_MAXMIN_*` flags are replaced with the `di_min_mapmem`, `di_des_mapmem`, and `di_max_mapmem` fields. They not only specify the required amount of DMA space, but also the necessary I/O memory entitlement for the device.

The `di_min_mapmem` parameter is the minimum amount of memory that the driver must be able to map for DMA in order to ensure the forward progress. The `d_map_init_ext` kernel service fails if the minimum I/O memory entitlement requirement cannot be satisfied.

The `di_des_mapmem` parameter is the required amount of memory that the driver wants to be able to I/O map in order to have good throughput. In most cases, this is a value that a driver specifies through the `DMA_MAXMIN_*` flag.

The `di_max_mapmem` parameter is the maximum amount of memory that the driver can ever map for DMA. This is the amount of DMA space that the `d_map_init_ext` kernel service can allocate.

**Note:** While the I/O memory entitlement for a `d_handle_t` type can be changed at runtime through the `d_map_attr` kernel service, the DMA space cannot be changed dynamically.

The `di_mapper_name` parameter contains the name of the device instance using the DMA resources (for example, `ent0`, `scsi1`, and so on).

## Execution Environment

The `d_map_init_ext` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful completion.
<code>struct d_handle *</code>	Indicates a successful completion.
<code>ENOMEM_D_MAP_INIT_EXT_1</code>	Indicates that the memory allocation failed. An AIX error is logged.
<code>ENOMEM_D_MAP_INIT_EXT_2</code>	Cannot reserve I/O memory entitlement with the least amount specified by the <code>di_min_mapmem</code> parameter. An AIX error is logged.
<code>ENOMEM_D_MAP_INIT_EXT_3</code>	Cannot allocate enough DMA space. An AIX error is logged.
<code>EINVAL_D_MAP_INIT_EXT_1</code>	Indicates that some input argument is not valid. An AIX error is logged in some cases.
<code>EINVAL_D_MAP_INIT_EXT_2</code>	Indicates that the combination of input arguments and system configuration is not valid. No AIX error is logged.
<code>EINVAL_D_MAP_INIT_EXT_3</code>	Indicates that the RAS initialization failed. No AIX error is logged.

### Related reference:

“`d_map_clear` Kernel Service” on page 80

“`d_unmap_page` Kernel Service” on page 108

“`d_map_list` Kernel Service”

## `d_map_list` Kernel Service

### Purpose

Performs platform-specific DMA mapping for a list of virtual addresses.

## Syntax

```
#include <sys/dma.h>

int d_map_list (*handle, flags, minxfer, *virt_list, *bus_list)
struct d_handle *handle;
int flags;
int minxfer;
struct dio *virt_list;
struct dio *bus_list;
```

**Note:** The following is the interface definition for `d_map_list` when the `DMA_ADDRESS_64` and `DMA_ENABLE_64` flags are set on the `d_map_init` call.

```
int d_map_list (*handle, flags, minxfer, *virt_list, *bus_list)
struct d_handle *handle;
int flags;
int minxfer;
struct dio_64 *virt_list;
struct dio_64 *bus_list;
```

## Parameters

Item	Description
<i>handle</i>	Indicates the unique handle returned by the <code>d_map_init</code> kernel service.
<i>flags</i>	Specifies one of the following flags:  <b>DMA_READ</b> Transfers from a device to memory.  <b>BUS_DMA</b> Transfers from one device to another device.  <b>DMA_BYPASS</b> Do not check page access.  <b>DMA_STMAP</b> Indicates a short-term mapping.
<i>minxfer</i>	Specifies the minimum transfer size for the device.
<i>virt_list</i>	Specifies a list of virtual buffer addresses and lengths.
<i>bus_list</i>	Specifies a list of bus addresses and lengths.

## Description

The `d_map_list` kernel service is a bus-specific utility routine determined by the `d_map_init` kernel service that accepts a list of virtual addresses and sizes and provides the resulting list of bus addresses. This service fills out the corresponding bus address list for use by the device in performing the DMA transfer. This service allows for scatter/gather capability of a device and also allows the device to combine multiple requests that are contiguous with respect to the device. The lists are passed via the `dio` structure. If the `d_map_list` service is unable to complete the mapping due to exhausting the capacity of the provided `dio` structure, the `DMA_DIOFULL` error is returned. If the `d_map_list` service is unable to complete the mapping due to exhausting resources required for the mapping, the `DMA_NORES` error is returned. In both of these cases, the `bytes_done` field of the `dio` virtual list is set to the number of bytes successfully mapped. This byte count is a multiple of the `minxfer` size for the device as provided on the call to `d_map_list`. The `resid_iov` field is set to the index of the remaining `d_iovec` fields in the list. Unless the `DMA_BYPASS` flag is set, this service verifies access permissions to each page. If an access violation is encountered on a page with the list, the `DMA_NOACC` error is returned, and the `bytes_done` field is set to the number of bytes preceding the faulting `iovec`. If the mapping is for short term (that is, it is unmapped as soon as the I/O is complete), you must set the `DMA_STMAP` flag.

### Note:

1. When the `DMA_NOACC` return value is received, no mapping is done, and the bus list is undefined. In this case, the `resid_iov` field is set to the index of the `d_iovec` that encountered the access violation.

2. You can use the `D_MAP_LIST` macro provided in the `/usr/include/sys/dma.h` file to code calls to the `d_map_list` kernel service.

## Return Values

Item	Description
<code>DMA_NORES</code>	Indicates that resources were exhausted during mapping.

**Note:** `d_map_list` possible partial transfer was mapped. Device driver may continue with partial transfer and submit the remainder on a subsequent `d_map_list` call, or call `d_unmap_list` to undo the partial mapping. If a partial transfer is issued, then the driver must call `d_unmap_list` when the I/O is complete.

Item	Description
<code>DMA_DIOFULL</code>	Indicates that the target bus list is full.

**Note:** `d_map_list` possible partial transfer was mapped. Device driver may continue with partial transfer and submit the remainder on a subsequent `d_map_list` call, or call `d_unmap_list` to undo the partial mapping. If a partial transfer is issued, then the driver must call `d_unmap_list` when the I/O is complete.

Item	Description
<code>DMA_NOACC</code>	Indicates no access permission to a page in the list.

**Note:** `d_map_list` no mapping was performed. No need for the device driver to call `d_unmap_list`, but the driver must fail the faulting I/O request, and resubmit any remainder in a subsequent `d_map_list` call.

Item	Description
<code>DMA_SUCC</code>	Indicates that the entire transfer successfully mapped.

**Note:** `d_map_list` successful mapping was performed. Device driver must call `d_unmap_list` when the I/O is complete. In the case of a long-term mapping, the driver must call `d_unmap_list` when the long-term mapping is no longer needed.

### Related reference:

“`d_map_init` Kernel Service” on page 82

“`d_map_init_ext` Kernel Service” on page 83

## `d_map_page` Kernel Service

### Purpose

Performs platform-specific DMA mapping for a single page.

### Syntax

```
#include <sys/dma.h>
#include <sys/xmem.h>

int d_map_page(*handle, flags, baddr, *busaddr, *xmp)
struct d_handle *handle;
int flags;
caddr_t baddr;
uint *busaddr;
struct xmem *xmp;
```

**Note:** The following is the interface definition for `d_map_page` when the `DMA_ADDRESS_64` and `DMA_ENABLE_64` flags are set on the `d_map_init` call.

```
int d_map_page(*handle, flags, baddr, *busaddr, *xmp)
struct d_handle *handle;
int flags;
unsigned long long baddr;
unsigned long long *busaddr;
struct xmem *xmp;
```

## Parameters

Item	Description
<i>handle</i>	Indicates the unique handle returned by the <code>d_map_init</code> kernel service.
<i>flags</i>	Specifies one of the following flags:  <b>DMA_READ</b> Transfers from a device to memory.  <b>BUS_DMA</b> Transfers from one device to another device.  <b>DMA_BYPASS</b> Do not check page access.  <b>DMA_STMAP</b> Indicates a short-term mapping.
<i>baddr</i>	Specifies the buffer address.
<i>busaddr</i>	Points to the <i>busaddr</i> field.
<i>xmp</i>	Cross-memory descriptor for the buffer.

## Description

The `d_map_page` kernel service is a bus-specific utility routine determined by the `d_map_init` or `d_map_init_ext` kernel service that performs platform specific mapping of a single 4KB or less transfer for DMA master devices. The `d_map_page` kernel service is a fast-path version of the `d_map_list` service. The entire transfer amount must fit within a single page in order to use this service. This service accepts a virtual address and completes the appropriate bus address for the device to use in the DMA transfer. Unless the `DMA_BYPASS` flag is set, this service also verifies access permissions to the page. If the mapping is for short term (that is, it is unmapped as soon as the I/O is complete), you must set the `DMA_STMAP` flag.

If the buffer is a global kernel space buffer, the cross-memory descriptor can be set to point to the exported `GLOBAL` cross-memory descriptor, *xmem\_global*.

If the transfer is unable to be mapped due to resource restrictions, the `d_map_page` service returns `DMA_NORES`. If the transfer is unable to be mapped due to page access violations, this service returns `DMA_NOACC`.

**Note:** You can use the `D_MAP_PAGE` macro provided in the `/usr/include/sys/dma.h` file to code calls to the `d_map_page` kernel service.

## Return Values

Item	Description
DMA_NORES	Indicates that resources are unavailable.

**Note:** `d_map_page` no mapping is done, device driver must wait until resources are freed and attempt the `d_map_page` call again.

Item	Description
DMA_NOACC	Indicates no access permission to the page.

**Note:** `d_map_page` no mapping is done, device driver must fail the corresponding I/O request.

Item	Description
DMA_SUCC	Indicates that the <code>busaddr</code> parameter contains the bus address to use for the device transfer.

**Note:** `d_map_page` successful mapping was done, device driver must call `d_unmap_page` when I/O is complete, or when device driver is finished with the mapped area in the case of a long-term mapping.

**Related reference:**

“`d_alloc_dmamem` Kernel Service” on page 61

“`d_map_init` Kernel Service” on page 82

“`d_map_list` Kernel Service” on page 84

## **d\_map\_query** Kernel Service

### **Purpose**

Queries the amount of direct memory access (DMA) space or DMA windows available on the partition end point. To use full 64-bit DMA all device drivers must call the `d_map_query` kernel service before attempting to initialize a new DMA window or before attempting to allocate DMA space within an existing DMA window by using the `d_map_init_ext` kernel service. Device drivers that do not use full 64-bit DMA should not call this service.

### **Syntax**

```
#include <sys/types.h>
#include <sys/dma.h>
#include <sys/kerrno.h>

kerrno_t d_map_query(bid, slot, flags, cmd, dq_info)
uint64_t bid;
uint64_t slot;
uint64_t flags;
uint64_t cmd;
void * dq_info;
```

### **Parameter**

Item	Description
<i>bid</i>	Specifies the bus identifier.
<i>slot</i>	Specifies the slot on the parent bus. This is the same as the <code>connwhere</code> property in the <code>CuDv</code> object class for the device.
<i>flags</i>	Specifies flags for the <code>d_map_query</code> kernel service. For future support, this parameter must be set to 0.
<i>cmd</i>	Specifies the type of query that the <code>d_map_query</code> kernel service will execute.
<i>dq_info</i>	Specifies the <code>dq_ddw_resources_t</code> or <code>dq_dma_available_t</code> structure based on which <code>cmd</code> parameter was defined.



## Description

The `d_map_query` kernel service allows the device driver to determine the amount of DMA space available within the DMA window or the amount of DMA windows available for a particular partition end point.

The `d_map_query` kernel service Dynamic DMA Windows Query (**DDW\_QUERY**) option is supported only on partition end points that support the dynamic DMA windows (**DDW\_QUERY**) option. The `d_map_query` kernel service can also be used to determine the dynamic DMA windows capability of a particular partition endpoint. When a slot is initialized on a reboot or power-on operation of the partition or on a DR isolate operation that encompasses the partition endpoint, a default DMA window is always allocated for less than 4 GB. After the first call to the `d_map_query` kernel service with the **DDW\_QUERY** option, the default DMA window is removed. This leaves no usable DMA window on the partition endpoint until the `d_map_init_ext` kernel service is called to initialize a new DMA window.

**Note:** The **DDW\_QUERY** option should only be used by device drivers that fully support the 64-bit DMA.

The caller of the `d_map_query` kernel service must pass the desired command to the `cmd` parameter and have the appropriate `dq_info` parameter initialized.

The options available for the `cmd` parameter are defined in the `<sys/dma.h>` header file, and are described as follows:

### DDW\_QUERY

Returns the number of additional DMA windows available for a partition endpoint. The `dq_ddw_resources_t` structure must be passed to the `dq_info` parameter for this command. The `dqdr_version` field in the structure should be assigned as **DQDR\_VERSION**.

### DMA\_QUERY

Returns the maximum amount of contiguous pages available for a given page size in all existing DMA windows. The `dq_dma_available_t` structure must be passed to the `dq_info` parameter for this command. The `dqda_version` field in the structure should be assigned as **DQDA\_VERSION** and the corresponding I/O page size for the query must be specified. The supported I/O page size for the DMA operation can be obtained from the `d_map_query` kernel service by running the **DDW\_QUERY** command.

The `dq_ddw_resources_t` and `dq_dma_available_t` structures are defined in the `<sys/dma.h>` as follows:

```
typedef struct dq_ddw_resources
{
    /* input by caller */
    uint64_t dqdr_version;
    /* returned to caller */
    uint64_t dqdr_supported_page_sizes;
    uint64_t dqdr_windows_avail;
    /* Amount of dynamic DMA windows available.
    * If DDW_QUERY, is not available
    * 0 will be returned.
    */
    uint64_t dqdr_max_pages; /* Largest number of contiguous pages available.*/
    uint64_t dqdr_rsvd1;      /* reserved for future use */
    uint64_t dqdr_rsvd2;      /* reserved for future use */
    uint64_t dqdr_rsvd3;      /* reserved for future use */
    uint64_t dqdr_rsvd4;      /* reserved for future use */
} dq_ddw_resources_t;
/*
 * The dq_dma_available structure is to be used in d_map_query with the
 * DMA_QUERY cmd specified
 */
typedef struct dq_dma_available
{
```

```

    /*input by caller */
    uint64_t dqda_version;
    uint64_t dqda_io_page_size;
    /* Page size in bytes, should only be equal to the supported pagesize */
    /* returned to caller for DMA_Query*/
    uint64_t dqda_pages_available;
    uint64_t dqda_rsvd1;        /* reserved for future use */
    uint64_t dqda_rsvd2;        /* reserved for future use */
    uint64_t dqda_rsvd3;        /* reserved for future use */
    uint64_t dqda_rsvd4;        /* reserved for future use */
} dq_dma_available_t;

```

## Execution Environment

The `d_map_query` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Success
Kernno	Error occurred

### Related reference:

“`d_map_clear` Kernel Service” on page 80

“`d_map_list` Kernel Service” on page 84

“`d_unmap_list` Kernel Service” on page 107

## `d_map_slave` Kernel Service

### Purpose

Accepts a list of virtual addresses and sizes and sets up the slave DMA controller.

### Syntax

```

#include <sys/dma.h>

int d_map_slave (*handle, flags, minxfer, *vlist, chan_flag)
struct d_handle *handle;
int flags;
int minxfer;
struct dio *vlist;
uint chan_flag;

```

### Parameters

Item	Description
<i>handle</i>	Indicates the unique handle returned by the <code>d_map_init</code> kernel service.
<i>flags</i>	Specifies one of the following flags: <ul style="list-style-type: none"> <li><b>DMA_READ</b> Transfers from a device to memory.</li> <li><b>BUS_DMA</b> Transfers from one device to another device.</li> <li><b>DMA_BYPASS</b> Do not check page access.</li> </ul>
<i>minxfer</i>	Specifies the minimum transfer size for the device.
<i>vlist</i>	Specifies a list of buffer addresses and lengths.
<i>chan_flag</i>	Specifies the device and bus specific flags for the transfer.

## Description

The **d\_map\_slave** kernel service accepts a list of virtual buffer addresses and sizes and sets up the slave DMA controller for the requested DMA transfer. This includes setting up the system address generation hardware for a specific slave channel to indicate the specified data buffers, and enabling the specific hardware channel. The **d\_map\_slave** kernel service is not an exported kernel service, but a bus-specific utility routine determined by the **d\_map\_init** kernel service and provided to the caller through the **d\_handle** structure.

This service allows for scatter/gather capability of the slave DMA controller and also allows the device driver to coalesce multiple requests that are contiguous with respect to the device. The list is passed with the **dio** structure. If the **d\_map\_slave** kernel service is unable to complete the mapping due to resource, an error, **DMA\_NORES** is returned, and the **bytes\_done** field of the **dio** list is set to the number of bytes that were successfully mapped. This byte count is guaranteed to be a multiple of the *minxfer* parameter size of the device as provided to **d\_map\_slave**. Also, the *resid\_iov* field is set to the index of the remaining *d\_iovec* that could not be mapped. Unless the **DMA\_BYPASS** flag is set, this service will verify access permissions to each page. If an access violation is encountered on a page within the list, an error, **DMA\_NOACC** is returned and no mapping is done. The *bytes\_done* field of the virtual list is set to the number of bytes preceding the faulting *iovec*. Also in this case, the *resid\_iov* field is set to the index of the *d\_iovec* entry that encountered the access violation.

The virtual addresses provided in the *vlist* parameter can be within multiple address spaces, distinguished by the cross-memory structure pointed to for each element of the **dio** list. Each cross-memory pointer can point to the same cross-memory descriptor for multiple buffers in the same address space, and for global space buffers, the pointers can be set to the address of the exported GLOBAL cross-memory descriptor, *xmem\_global*.

The *minxfer* parameter specifies the absolute minimum data transfer supported by the device( the device blocking factor). If the device supports a minimum transfer of 512 bytes (floppy and disks, for example), the *minxfer* parameter would be set to 512. This allows the underlying services to map partial transfers to a correct multiple of the device block size.

### Note:

1. The **d\_map\_slave** kernel service does not support more than one outstanding DMA transfer per channel. Attempts to do multiple slave mappings on a single channel will corrupt the previous mappings.
2. You can use the **D\_MAP\_SLAVE** macro provided in the `/usr/include/sys/dma.h` file to code calls to the **d\_map\_clear** kernel service.
3. The possible flag values for the *chan\_flag* parameter can be found in `/usr/include/sys/dma.h`. These flags can be logically ORed together to reflect the desired characteristics of the device and channel.
4. If the **CH\_AUTOINIT** flag is used then the transfer described by the *vlist* pointer is limited to a single buffer address with a length no greater than 4K bytes.

## Return Values

Item	Description
DMA_NORES	Indicates that resources were exhausted during the mapping.
DMA_NOACC	Indicates no access permission to a page in the list.
DMA_BAD_MODE	Indicates that the mode specified by the <i>chan_flag</i> parameter is not supported.

#### Related reference:

“d\_map\_init Kernel Service” on page 82

## dmp\_add Kernel Service

### Purpose

Specifies data to be included in a system dump by adding an entry to the master dump table. Callers should use the “dmp\_ctl Kernel Service” on page 96. This service is provided for compatibility purposes.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/dump.h>
```

```
int dmp_add
( cdt_func)
struct cdt * ( (*cdt_func) ( ));
```

### Description

Kernel extensions use the **dmp\_add** service to register data areas to be included in a system dump. The **dmp\_add** service adds an entry to the master dump table. A master dump table entry is a pointer to a function provided by the kernel extension that will be called by the kernel dump routine when a system dump occurs. The function must return a pointer to a component dump table structure.

When a dump occurs, the kernel dump routine calls the function specified by the *cdt\_func* parameter twice. On the first call, an argument of 1 indicates that the kernel dump routine is starting to dump the data specified by the component dump table. On the second call, an argument of 2 indicates that the kernel dump routine has finished dumping the data specified by the component dump table. Kernel extensions should allocate and pin their component dump tables and call the **dmp\_add** service during initialization. The entries in the component dump table can be filled in later. The *cdt\_func* routine must not attempt to allocate memory when it is called.

**Note:** In AIX Version 7.1, this function automatically serializes CDT functions with I/O during dump time. The need for this function is device specific. Only the developer of the device can determine if this routine needs to be used. It is only recommended for devices that can be on the dump I/O path. Serializing I/O during dump time might degrade dump performance. Devices that are not on the dump path must either use the **dmp\_ctl** routine or the RAS system dump interface.

### The Component Dump Table

The component dump table structure specifies memory areas to be included in the system dump. The structure type (**struct cdt**) is defined in the */usr/include/sys/dump.h* file. A **cdt** structure consists of a fixed-length header (**cdt\_head** structure) and an array of one or more **cdt\_entry** structures. The **cdt\_head** structure contains a component name field, which should be filled in with the name of the kernel extension, and the length of the component dump table. Each **cdt\_entry** structure describes a contiguous data area, giving a pointer to the data area, its length, a segment register, and a name for the data area.

### Use of the Formatting Routine

Each kernel extension that includes data in the system dump can install a unique formatting routine in the `/var/adm/ras/dmprtns` directory. The name of the formatting routine must match the component name field of the corresponding component dump table.

The dump image file includes a copy of each component dump table used to dump memory. A sample dump formatter is shipped with `bos.sysmgt.serv_aid` in the `/usr/samples/dumpfmt` directory.

### Organization of the Dump Image File

Memory dumped for each kernel extension is laid out as follows in the dump image file. The component dump table is followed by a bitmap for the first data area, then the first data area itself, then a bitmap for the next data area, the next data area itself, and so on.

The bitmap for a specific data area indicates which pages of the data area are present in the dump image and which are not. Pages that were not in memory when the dump occurred were not dumped. The least significant bit of the first byte of the bitmap is set to 1 (one) if the first page is present. The next least significant bit indicates the presence or absence of the second page and so on.

A macro for determining the size of a bitmap is provided in the `/usr/include/sys/dump.h` file.

### Parameters

Item	Description
<code>cdt_func</code>	Specifies a function that returns a pointer to a component dump table entry. The function and the component dump table entry both must reside in pinned global memory.

### Execution Environment

The `dmp_add` kernel service can be called from the process environment only.

### Return Values

Item	Description
0	Indicates a successful operation.
-1	Indicates that the function pointer to be added is already present in the master dump table.

#### Related reference:

“`dmp_ctl` Kernel Service” on page 96

#### Related information:

`exec`: `execl`, `execle`, `execlp`, `execv`, `execve`, `execvp`, or `exec` Subroutine

RAS Kernel Services

## `dmp_compspec` and `dmp_compext` Kernel Services

### Purpose

Specifies a component and callback parameters to be included in the dump.

### Syntax

```
#include <sys/livedump.h>
```

```
kernno_t dmp_compspec (flags, comp, anchor, extid, p1, p2, ..., NULL)
long flags;
long comp;
void *anchor;
dmp_extid_t *extid;
```

```
char *p1;
char *p2;
...
```

```
kernno_t dmp_compext (extid, p1, p2, ..., NULL)
dmp_extid_t extid;
char *p1;
char *p2;
...
```

## Parameters

Item	Description
<i>anchor</i>	Points to the associated <b>ldmp_parms_t</b> data structure or to an <b>ldmp_prepare_t</b> data structure.
<i>comp</i>	Specifies the component, specified as indicated by the flags.
<i>extid</i>	Points to an item of <b>dmp_extid_t</b> type, for the <b>dmp_compspec</b> kernel service, where an identifier is returned, if you use the <b>dmp_compext</b> kernel service to provide additional parameters for the component being dumped. This identifier might then be specified to add additional parameters to the component using the <b>dmp_compext</b> kernel service. The <i>extid</i> parameter can be NULL.
<i>flags</i>	You can specify the following values: <p><b>DCF_FAILING</b> Indicates that this is the failing component. You can only specify one failing component.</p> <p><b>DCF_FIRST</b> Indicates that this component is to be dumped first. Normally components are dumped in the order specified. <b>Note:</b></p> <ul style="list-style-type: none"> <li>• The <b>DCF_FIRST</b> value is only valid when the anchor refers to an <b>ldmp_parms_t</b> data item. It is not valid when the callback receives the <b>RASCD_LDMP_PREPARE</b> command.</li> <li>• The last component specified to be dumped first is the one dumped first.</li> </ul> <p><b>DCF_LEVEL0 - DCF_LEVEL9</b> Indicates the detail level, 0 through 9, to dump this component. If none of these flags are set, the component is dumped at its current level.</p> <p><b>DCF_MINIMAL</b> Indicates the DCF_LEVEL1 level.</p> <p><b>DCF_NORMAL</b> Indicates the DCF_LEVEL3 level.</p> <p><b>DCF_DETAIL</b> Indicates the DCF_LEVEL7 level.</p> <p><b>DCF_LONG</b> Indicates that the parameters are two parameters of long type. Rather than passing in an unlimited number of strings, a component can be passed in two long data items, as in the case with pseudo-components. One and only one of the following component specification flags must be given. They specify how the component is specified in the <i>dc_component</i> field:</p> <p><b>DCF_BYNAME</b> Indicates that the component is specified by path name.</p> <p><b>DCF_BYLNAME</b> Indicates that the component is specified by logical alias.</p> <p><b>DCF_BYTYPE</b> Indicates that the component is specified by type.</p> <p><b>DCF_BYCB</b> Indicates that the component is specified by <i>ras_block_t</i>.</p>

<b>Item</b>	<b>Description</b>
<i>p1</i> , <i>p2</i> ...	Specifies the component's parameters, the last of which must be NULL. If keyword parameters are being specified, The parameters must be strings, and contain the keyword and its values. If multiple keyword and value pairs appear in a single parameter, they are separated with blanks. For example, the <i>p1</i> parameter can be foo=1234, and the <i>p2</i> parameter can be bar=5678,16. Also, the <i>p1</i> parameter can be foo=1234 bar=5678.
	If the <b>DCF_LONG</b> flag is set, two parameters of long type are passed in. In this case, the <i>p1</i> and <i>p2</i> parameters contain the values of long type, and no more parameters can be specified.

## Description

The **dmp\_compspec** and **dmp\_compevt** kernel services provide components and their callback parameters for a dump. You can only use these kernel services in a live dump.

The **dmp\_compspec** kernel service is used before you start a live dump with the **livedump** kernel service. You can also use this kernel service when a component's callback wants to include another component in a live dump, that is, when the callback receives the **RASCD\_LDMP\_PREPARE** command.

Multiple components can be included in a live dump.

The **dmp\_compevt** function is used to provide additional parameters for a component.

## Return Values

<b>Item</b>	<b>Description</b>
0	Indicates a successful completion.
EINVAL_RAS_DMP_COMPSPEC_FLAGS	Indicates that the flags specification is not valid.
EINVAL_RAS_DMP_COMPSPEC_COMP	Indicates that the component specification is not valid.
EINVAL_RAS_DMP_COMPSPEC_NOTAWARE	Indicates that the specified component must support live dump.
EINVAL_RAS_DMP_COMPSPEC_ANCHOR	Indicates that the anchor specification is not valid.
EFAULT_RAS_DMP_COMPSPEC_ANCHOR	Indicates that the storage the anchor parameter refers to is not valid.
EFAULT_RAS_DMP_COMPSPEC_EXTID	Indicates that the storage the <i>extid</i> parameter refers to is not valid.
EFAULT_RAS_DMP_COMPSPEC_PARMS	Indicates that a parameter address is not valid.
EINVAL_RAS_LDMP_ESTIMATE	Indicates that the anchor parameter indicates a dump size estimate request, but the <b>dmp_compspec</b> call was not made from the process environment.
EINVAL_RAS_DMP_COMPSPEC_NOADD	Indicates that components cannot be added to this dump, that is, the dump type flags, <i>ldpr_flags</i> , have the <b>LDT_NOADD</b> bit set.
EINVAL_RAS_DMP_COMPSPEC_FAILING	Indicates that the failing component has already been specified.
ENOMEM_RAS_DMP_COMPSPEC	Indicates that no storage is available.
EINVAL_RAS_DMP_COMPEXT_EXTID	Indicates that the <i>extid</i> parameter does not refer to a valid component.
EFAULT_RAS_DMP_COMPEXT_EXTID	Indicates that the storage the <i>extid</i> parameter refers to is not valid.
EFAULT_RAS_DMP_COMPEXT_PARMS	Indicates that the storage a parameter refers to is not valid.
EBUSY_RAS_DMP_COMPEXT	Indicates that the specification of this component is complete, and no more parameters can be added. This happens if the component the <i>extid</i> parameter referred to has already completed its <b>RASCD_LDMP_PREPARE</b> processing.
ENOMEM_RAS_DMP_COMPEXT	Indicates that no storage is available.

### Related reference:

“livedump Kernel Service” on page 337

“ldmp\_setupparms Kernel Service” on page 335

“ras\_ret\_query\_parms Kernel Service” on page 434

## dmp\_ctl Kernel Service

### Purpose

Adds and removes entries to the master dump table.

### Syntax

```
#include <sys/types.h>
#include <errno.h>
#include <sys/dump.h>

int dmp_ctl(op, parmp)
int op;
struct dmpctl_data *parmp;
```

### Description

The **dmp\_ctl** kernel service is used to manage dump routines. It replaces the **dmp\_add** and **dmp\_del** kernel services which are still supported for compatibility reasons. The major differences between routines added with the **dmp\_add()** command and those added with the **dmp\_ctl()** command are:

- The routines are invoked differently from routines added with the **dmp\_add** kernel service. Routines added using the **dmp\_ctl** kernel service return a void pointer, to a dump table or to a dump size estimate.
- Routines added with the **dmp\_ctl** kernel service are expected to ignore functions they don't support. For example, they should not trap if they receive an unrecognized request. This allows future functionality to be added without all users needing to change.

The **dmp\_ctl** kernel service is used to request that an amount of memory be set aside in a global buffer. This will then be used by the routine to store data not resident in memory. An example of such data is dump data provided by an adapter. Without a global buffer, the data would need to be placed into a pinned buffer allocated at configuration time. Each component would need to allocate its own pinned buffer.

The system dump facility maintains a global buffer for such data. This buffer is allocated when it is first requested, with the requested size. Another dump routine requesting more data causes the buffer to be reallocated with the larger size. Since this buffer must be maintained in pinned storage for the life of the system, only ask for as much memory as is required. Asking for an excessive amount of storage will compromise system performance by reserving too much pinned storage.

Any dump routine using the global buffer is called whenever dump data is required. Routines are only called once to provide such data. Their dump table addresses are saved and used if the dump is restarted.

**Note:** The **dmp\_ctl** kernel service can also be used by a dump routine to report a routine failure. This may be necessary if the routine detects that it can't dump what needs to be dumped for some reason such as corruption of a data structure.

**Note:** Beginning with AIX Version 6.1 with the 6100-02 Technology Level, the **dmp\_ctl** kernel service supports that DMPFUNC\_SERIALIO operation flag.

### Dump Tables

A dump routine returns a component dump table that begins with **DMP\_MAGIC**, which is the magic number for the 32- or 64-bit dump table. If the unlimited sized dump table is used, the magic number is **DMP\_MAGIC\_U** and the **cdt\_u** structure is used. If this is the case, the dump routine is called repeatedly until it returns a null **cdt\_u** pointer. The purpose of the unlimited size dump table is to provide a way to dump an unknown number of data areas without having to preallocate the largest possible array of



`cdt_entry` elements as is required for the classic dump table. The definitions for dump tables are in the `sys/dump.h` include file.

## Parameters

`dmp_ctl` operations and the `dmpctl_data` structure are defined in the `dump.h` text file.

Item	Description
op	Specifies the operation to perform.
parmp	Points to a <code>dmpctl_data</code> structure containing values for the specified operation. The <code>dmpctl_data</code> structure is defined in the <code>/usr/include/sys/dump.h</code> file as follows: <pre> /* Dump Routine failures data. */ struct __rtnf {     int rv;                /* error code. */     ulong vaddr;          /* address. */     vmhandle_t handle;    /* handle */ };  typedef void *(*__CDTFUNCENH)(int op, void *buf); struct dmpctl_data {     int dmpc_magic;        /* magic number */     int dmpc_flags;        /* dump routine flags. */     __CDTFUNCENH dmpc_func;     union {         u_longlong_t bsize; /* Global buffer size requested. */         struct __rtnf rtnf;     } dmpc_u; }; #define DMPC_MAGIC1 0xcdcdc01 #define DMPC_MAGIC DMPC_MAGIC1 #define dmpc_bsize dmpc_u.bsize #define dmpcf_rv dmpc_u.rtnf.rv #define dmpcf_vaddr dmpc_u.rtnf.vaddr #define dmpcf_handle dmpc_u.rtnf.handle </pre>

The supported operations and their associated data are:

Item	Description
<b>DMPCTL_ADD</b>	<p>Adds the specified dump routine to the master dump table. This requires a pointer to the function and function type flags. Supported type flags are:</p> <p><b>DMPFUNC_CALL_ON_RESTART</b> Calls this function again if the dump is restarted. A dump function is only called once to provide dump data. If the function must be called and the dump is restarted on the secondary dump device, then this flag must be set. The <b>DMPFUNC_CALL_ON_RESTART</b> flag must be set if this function uses the global dump buffer. It also must be set if the function uses an unlimited size dump table, a table with <code>DMP_MAGIC_U</code> as the magic number.</p> <p><b>DMPFUNC_GLOBAL_BUFFER</b> This function uses the global dump buffer. The size is specified using the <code>dmpc_bsize</code> field.</p> <p><b>DMPFUNC_SERIALIO</b> Enables serialized I/O during dump time. The need for this flag is device specific. Only the developer of the device can determine if this flag needs to be set. It is only recommended for devices that can be on the dump I/O path. Serializing I/O during dump time can degrade dump performance. The default, without this flag, is to allow I/O to occur in parallel with CDT function calls.</p>
<b>DMPCTL_DEL</b>	Deletes the specified dump function from the master dump table.
<b>DMPCTL_RTNFAILURE</b>	Reports an inability to dump required data. The routine must set the <code>dmpc_func</code> , <code>dmpcf_rv</code> , <code>dmpcf_vaddr</code> , and <code>dmpcf_handle</code> fields.

Dump function invocation parameters:

Item	Description
operation code	<p>Specifies the operation the routine is to perform. Operation codes are:</p> <p><b>DMPRTN_START</b> The dump is starting for this dump table. Provide data.</p> <p><b>DMPRTN_DONE</b> The dump is finished. This call is provided so that a dump routine can do any cleanup required after a dump. This is specific to a device for which information was gathered. It does not free memory, since such memory must be allocated before the dump is taken.</p> <p><b>DMPRTN_AGAIN</b> Provide more data for this unlimited dump table. The routine must have first passed back a dump table beginning with DMP_MAGIC_U. When finished, the function must return a NULL.</p> <p><b>DMPRTN_ESTIMATE</b> Provide a size estimate. The function must return a pointer to an item of type dmp_sizeest_t. See the examples later in this article.</p>
buffer pointer	This is a pointer to the global buffer, or NULL if no global buffer space was requested.

## Return Values

Item	Description
0	Returned if successful.
EINVAL	Returned if one or more parameter values are invalid.
ENOMEM	Returned if the global buffer request can't be satisfied.
EEXIST	Returned if the dump function has already been added.

## Examples

- To add a dump routine (dmp\_rtn) that can be called once to provide data, type:

```

void *dmp_rtn(int op, void *buf);
struct cdt cdt;
dmp_sizeest_t estimate;

config()
{
    struct dmpctl_data parm;
    ...

    parm.dmpc_magic = DMPC_MAGIC1;
    parm.dmpc_func = dmp_rtn;
    parm.dmpc_flags = 0;
    ret = dmp_ctl(DMPCTL_ADD, &parm);
    ...
}

/*
 * Dump routine.
 *
 * input:
 *   op - dump routine operation.
 *   buf - NULL since no global buffer is used.
 *
 * returns:
 *   A pointer to the component dump table.
 */
void *
dmp_rtn(int op, void *buf)
{
    void *ret;

```

```

switch(op) {
case DMPRTN_START: /* Provide dump data. */
    ...
    ret = (void *)&cdt;
    break;
case DMPRTN_ESTIMATE:
    ret = (void *)&estimate;
    break;
default:
    break;
}

return(ret);
}

```

2. To add a dump routine (dmp\_rtn) that requests 16 kb of global buffer space, type:

```

...
#define BSIZ 16*1024
dmp_sizeest_t estimate;

config()
{
    ...
    parm.dmpc_magic = DMPC_MAGIC1;
    parm.dmpc_func = dmp_rtn;
    parm.dmpc_flags = DMPFUNC_CALL_ON_RESTART|DMPC_GLOBAL_BUFFER;
    parm.dmpc_bsize = BSIZ;
    ret = dmp_ctl(DMPCTL_ADD, &parm);
    ...
}

/*
 * Dump routine.
 *
 * input:
 * op - dump routine operation.
 * buf - points to the global buffer.
 *
 * output:
 * Return a pointer to the dump table or to the estimate.
 */
void *
dmp_rtn(int op, void *buf)
{
    void *ret;

    switch(op) {
case DMPRTN_START: /* Provide dump data. */
    ...
    (Put data in buffer at buf.)
    ret = (void *)&cdt;
    break;
case DMPRTN_ESTIMATE:
    ret = (void *)&estimate;
    break;
default:
    break;
}

return(ret);
}

```

**Related reference:**

“dmp\_del Kernel Service” on page 100

**Related information:**

Dump Special File  
System Dump Facility

## **dmp\_del Kernel Service**

### **Purpose**

Deletes an entry from the master dump table. Callers should use the “dmp\_ctl Kernel Service” on page 96. This service is provided for compatibility purposes.

### **Syntax**

```
#include <sys/types.h>  
#include <sys/errno.h>  
#include <sys/dump.h>
```

```
dmp_del ( cdt_func_ptr )  
struct cdt * ( (*cdt_func_ptr) ( ) );
```

### **Description**

Kernel extensions use the **dmp\_del** kernel service to unregister data areas previously registered for inclusion in a system dump. A kernel extension that uses the “dmp\_add Kernel Service” on page 92 to register such a data area can use the **dmp\_del** service to remove this entry from the master dump table.

### **Parameters**

Item	Description
<i>cdt_func_ptr</i>	Specifies a function that returns a pointer to a component dump table. The function and the component dump table must both reside in pinned global memory.

### **Execution Environment**

The **dmp\_del** kernel service can be called from the process environment only.

### **Return Values**

Item	Description
0	Indicates a successful operation.
-1	Indicates that the function pointer to be deleted is not in the master dump table.

### **Related reference:**

“dmp\_add Kernel Service” on page 92

“dmp\_ctl Kernel Service” on page 96

### **Related information:**

RAS Kernel Services

## **dmp\_eaddr, dmp\_context, dmp\_tid, dmp\_pid, dmp\_errbuf, dmp\_mtrc, dmp\_systrace, and dmp\_ct Kernel Services**

### **Purpose**

Provides functions for common dump tasks.

### **Syntax**

```
#include <sys/dump.h>
```

```
kernno_t dmp_eaddr (flags, anchor, name, addr, sz)
long flags;
void *anchor;
char *name;
long addr;
long sz;
```

```
kernno_t dmp_context (flags, anchor, name, ctx_type, p2)
long flags;
void *anchor;
char *name;
long ctx_type;
long p2;
```

```
kernno_t dmp_tid (flags, anchor, name, tid, unused)
long flags;
void *anchor;
char *name;
tid_t tid;
void *unused;
```

```
kernno_t dmp_pid (flags, anchor, name, pid, unused)
long flags;
void *anchor;
char *name;
pid_t pid;
void *unused;
```

```
kernno_t dmp_errbuf (flags, anchor, name, erridx, unused)
long flags;
void *anchor;
char *name;
ulong erridx;
long unused;
```

```
kernno_t dmp_mtrc (flags, anchor, name, com_sz, rare_sz)
long flags;
void *anchor;
char *name;
size_t com_sz;
size_t rare_sz;
```

```
kernno_t dmp_systrace (flags, anchor, name, sz, unused)
long flags;
void *anchor;
char *name;
long sz;
long unused;
```

```
kernno_t dmp_ct (flags, anchor, name, rasb, sz)
long flags;
void *anchor;
char *name;
ras_block_t rasb;
size_t sz;
```

## Parameters

Item	Description
<i>anchor</i>	Points to the associated <b>ldmp_parms_t</b> data structure or to an <b>ldmp_prepare_t</b> data structure.
<i>flags</i>	The <i>flags</i> parameter can be one or more of the following values: <b>DCF_FIRST</b> Specifies that this component is to be dumped first. Normally components are dumped in the order specified. <b>Note:</b> The last component specified to be dumped first is the one dumped first. <b>DCF_LEVEL0 - DCF_LEVEL9</b> Dumps the component at the specified detail level, 0 through 9. If none of these flags are set, the component is dumped at CD_LVL_NORMAL, detail level 3.
<i>name</i>	Specifies the name of the pseudo-component's dump table in the dump. The <i>name</i> parameter is only valid for the <b>dmp_eaddr</b> kernel service. You must specify the <i>name</i> parameter to NULL for the <b>dmp_context</b> , <b>dmp_fid</b> , <b>dmp_pid</b> , <b>dmp_errbuf</b> , <b>dmp_mtrc</b> , <b>dmp_systrace</b> , and <b>dmp_ct</b> kernel services.
<i>unused</i>	You must specify this parameter to NULL or 0.
The remaining parameters are pseudo-component dependent: <b>dmp_eaddraddr</b>	Specifies the effective address of the memory to be dumped.
<i>sz</i>	Specifies the length of the memory in bytes.
<b>dmp_contextctx_type</b>	Specifies the context to dump. It can be one of the following values: <b>DMP_CTX_CUR</b> To dump the current context. <b>DMP_CTX_PREV</b> To dump the previous context. <b>DMP_CTX_SPEC</b> To dump the context specified by the <i>p2</i> parameter. The <i>p2</i> parameter must contain the address of the <b>ksmtsava</b> structure for the context. <b>DMP_CTX_RWA</b> To dump the context from the supplied recovery work area. The <i>p2</i> parameter must contain the address of the recovery work area, <i>rwa</i> . <b>DMP_CTX_BID or DMP_CTX_LCPUID</b> To dump the context for the processor specified by the <i>p2</i> parameter. You can specify the processor either by the bind ID or by the logical ID. <b>DMP_CTX_TID</b> To dump the context of the thread specified by the <i>p2</i> parameter, which must contain the thread ID.
<i>p2</i>	Specifies the address of the context, the logical processor ID, the bind ID, or the thread ID dependent on the value of the <i>ctx_type</i> parameter.
<b>dmp_tidtid</b>	Specifies the ID of the thread to dump.
<b>dmp_pidpid</b>	Specifies the ID of the process to dump.
<b>dmp_errbuferridx</b>	Specifies the kernel workload partition (WPAR) ID of the partition's error logging buffer to dump. The value of 0 stands for the global buffer.
<b>dmp_mtrccom_sz</b>	Specifies the amount of common to dump.
<i>rare_sz</i>	Specifies the amount of rare data to dump.
<b>dmp_systracesz</b>	Specifies the amount of system trace data to dump. If the <i>sz</i> parameter is set to 0, all the buffered trace data is dumped, up to the amount allowed by the detail level.
<b>dmp_ctrasb</b>	Specifies the <i>ras_block_t</i> of the component whose component trace is to be dumped.
<i>sz</i>	Specifies the amount of data to dump. If the <i>sz</i> parameter is set to 0, all the components' trace data is dumped, up to the limit for the detail level.

## Description

The **dmp\_eaddr** kernel service dumps memory by effective address.

The **dmp\_context** kernel service dumps the specified thread context.

The **dmp\_tid** kernel service dumps the kernel data for a thread.

The **dmp\_pid** kernel service dumps the kernel data for a process.

The **dmp\_errbuf** kernel service dumps the error logging buffer for the specified partition.

The **dmp\_mtrc** kernel service dumps entries from the lightweight memory trace buffers.

The **dmp\_systrace** dumps entries from the system trace buffers.

The **dmp\_ct** dumps component trace entries.

## Execution Environment

The **dmp\_eaddr**, **dmp\_context**, **dmp\_tid**, **dmp\_pid**, **dmp\_errbuf**, **dmp\_mtrc**, **dmp\_systrace**, and **dmp\_ct** kernel services can be called from either the process or interrupt environment.

## Return Values

Item	Description
0	Indicates a successful completion.
EINVAL_DMP_PSEUDO	Indicates that the name parameter is not valid.
EINVAL_DMP_CHECK_ANCHOR	Indicates that no anchor was specified, or the anchor parameter does not point to an area of <b>ldmp_parms_t</b> or <b>ldmp_prepare_t</b> type.
EFAULT_DMP_CHECK_ANCHOR	Indicates that the storage specified by the <i>anchor</i> parameter is not valid.
EINVAL_RAS_DMP_COMPSPEC_FLAGS	Indicates that the flags specification is not valid. This error also occurs if the <b>DCF_FIRST</b> flag is specified when the anchor is an <b>ldmp_prepare_t</b> data item.
EINVAL_RAS_DMP_COMPSPEC_NOADD	Indicates that components cannot be added to this dump.
ENOMEM_RAS_DMP_COMPSPEC	Indicates that the storage is not sufficient.
EINVAL_RAS_DMP_EADDR	Indicates that the flags parameter is not valid.
EINVAL_RAS_DMP_CONTEXT	Indicates that the parameter of the <b>dmp_context</b> kernel service is not valid. This is also returned if the <i>p2</i> parameter is not used, but is not NULL.
ENOENT_RAS_DMP_CONTEXT_CTX_NOTFOUND	Indicates that the specified context was not found.
EFAULT_RAS_DMP_CONTEXT	Indicates that the storage the specified context pointer points to is not valid.
EINVAL_RAS_DMP_TID	Indicates that the parameter of the <b>dmp_tid</b> kernel service is not valid.
EINVAL_RAS_DMP_PID	Indicates that the parameter of the <b>dmp_pid</b> kernel service is not valid.
EINVAL_RAS_DMP_ERRBUF	Indicates that the parameter of the <b>dmp_errbuf</b> kernel service is not valid.
ECHRNG_RAS_DMP_ERRBUF	Indicates that the <i>erridx</i> parameter is out of range.
EINVAL_RAS_DMP_MTRC	Indicates that the parameter of the <b>dmp_mtrc</b> kernel service is not valid.
ENOENT_RAS_DMP_MTRC	Indicates that the lightweight memory trace is not active.
EINVAL_RAS_DMP_SYSTRACE	Indicates that the parameter of the <b>dmp_systrace</b> kernel service is not valid.
ENOENT_RAS_DMP_SYSTRACE	Indicates that the system trace is not active.
EINVAL_RAS_DMP_CT	Indicates that the parameter of the <b>dmp_ct</b> kernel service is not valid.
ENOMEM_RAS_DMP_CT	Indicates that the storage is not sufficient.
EINVAL_RAS_DMP_CT_GETPATH	Indicates that the specified component is not valid.
EINVAL_RAS_DMP_CT_LOOKUP	Indicates that an error occurred while this component was being validated.
ENOTSUP_RAS_DMP_CT	Indicates that the specified component does not have a component trace.

## Related Information

The **livedump** kernel service and **dmp\_kernext** kernel service.

### Related reference:

“livedump Kernel Service” on page 337

“dmp\_kernext Kernel Service” on page 104

## dmp\_kernext Kernel Service

### Purpose

Causes the specified kernel extension to be shipped with the live dump for symbol resolution.

### Syntax

```
#include <sys/dump.h>
```

```
kernno_t dmp_kernext (anchor, ptr)void *anchor;  
void *ptr;
```

### Parameters

Item	Description
<i>anchor</i>	Points to either an <code>ldmp_parms_t</code> or <code>ldmp_prepare_t</code> structure.
<i>ptr</i>	Specifies an address within the kernel extension. If the value is 0, the dump includes information for all loaded kernel extensions.

### Description

The `dmp_kernext` kernel service causes `snap` to package the specified kernel extension with the current live dump. This also includes loader information for the extension in the dump. You can specify the extension by setting the `ptr` parameter to a text or data address within the extension. The extension's file name is noted in the dump, and `snap` can be used to cause this file to be bundled with the `snap` data when the dump is collected for sending to IBM®.

### Execution Environment

The `dmp_kernext` kernel service can be called from either the process or interrupt environment.

### Return Values

Item	Description
0	Indicates a successful completion.
EINVAL_RAS_DMP_KERNEXT	Indicates that the <i>anchor</i> parameter is not valid.

### Related reference:

“livedump Kernel Service” on page 337

### Related information:

`snap` subroutine

## d\_roundup Kernel Service

### Purpose

Rounds the value `length` up to a given number of cache lines.

### Syntax

```
int d_roundup(length)  
int length;
```

### Parameter



Item	Description
<i>length</i>	Specifies the size in bytes to be rounded.

## Description

To maintain cache consistency, buffers must occupy entire cache lines. The **d\_roundup** service helps provide that function by rounding the value *length* up to a given number in integer form.

## Execution Environment

The **d\_roundup** service can be called from either the process or interrupt environment.

### Related reference:

“d\_align Kernel Service” on page 61

“d\_cflush Kernel Service” on page 62

### Related information:

Understanding Direct Memory Access (DMA) Transfers

## d\_sync\_mem Kernel Service

### Purpose

Allows a device driver to indicate that previously mapped buffers may need to be refreshed.

### Syntax

```
int d_sync_mem(d_handle_t handle, dio_t blist)
```

### Description

The **d\_sync\_mem** service allows a device driver to indicate that previously mapped buffers may need to be refreshed, either because a new DMA is about to start or a previous DMA has now completed. **d\_sync\_mem** is not an exported kernel service, but a bus-specific utility determined by **d\_map\_init** based on platform characteristics and provided to the caller through the *d\_handle* structure. **d\_sync\_mem** allows the driver to identify additional coherency points beyond those of the initial mapping (**d\_map\_list**) and termination of the mapping (**d\_unmap\_list**). Thus **d\_sync\_mem** provides a way to long-term map buffers and still handle potential data consistency problems.

The *blist* parameter is a pointer to the **dio** structure that describes the initial mapping, as returned by **d\_map\_list**. Note that for bounce buffering, the data direction is also implicitly defined by this initial mapping.

- If the **map\_list** call describes a transfer from system memory to a device, subsequent **d\_sync\_mem** calls using the corresponding *blist* will synchronize the memory view. This assumes that the original system memory pages contain the correct data.
- If the **map\_list** call describes a transfer from a device to system memory, then subsequent **d\_sync\_mem** calls will synchronize the memory view. This assumes that the bounce pages the device directly accessed contained the correct data.

**Note:** You can use the **D\_SYNC\_MEM** macro provided in the `/usr/include/sys/dma.h` file to code calls to the **d\_sync\_mem** kernel service.

### Parameters

Item	Description
<i>d_handle_t</i>	Indicates the unique dma handle returned by <code>d_map_init</code> .
<i>dio_t blist</i>	List of vectors returned by original <code>d_map_list</code> .

## Return Values

Item	Description
<code>DMA_SUCC</code>	Buffers described by the <i>blist</i> have been synchronized.
<code>DMA_FAIL</code>	Buffers could not be synchronized.

### Related reference:

“`d_alloc_dmamem` Kernel Service” on page 61

“`d_map_list` Kernel Service” on page 84

“`d_unmap_list` Kernel Service” on page 107

## DTOM Macro for mbuf Kernel Services

### Purpose

Converts an address anywhere within an `mbuf` structure to the head of that `mbuf` structure.

### Syntax

```
#include <sys/mbuf.h>
```

```
DTOM ( bp );
```

### Parameter

Item	Description
<i>bp</i>	Points to an address within an <code>mbuf</code> structure.

### Description

The `DTOM` macro converts an address anywhere within an `mbuf` structure to the head of that `mbuf` structure. This macro is valid only for `mbuf` structures without an external buffer (that is, with the `M_EXT` flag not set).

This macro can be viewed as the opposite of the `MTOD` macro, which converts the address of an `mbuf` structure into the address of the actual data contained in the buffer. However, the `DTOM` macro is more general than this view implies. That is, the input parameter can point to any address within the `mbuf` structure, not merely the address of the actual data.

### Example

The `DTOM` macro can be used as follows:

```
char          *bp;
struct mbuf   *m;
m = DTOM(bp);
```

### Related reference:

“`MTOD` Macro for mbuf Kernel Services” on page 376

### Related information:

I/O Kernel Services

## **d\_unmap\_list Kernel Service**

### **Purpose**

Deallocates resources previously allocated on a **d\_map\_list** call.

### **Syntax**

```
#include <sys/dma.h>
```

```
void d_unmap_list (*handle, *bus_list)
struct d_handle *handle
struct dio *bus_list
```

**Note:** The following is the interface definition for **d\_unmap\_list** when the **DMA\_ADDRESS\_64** and **DMA\_ENABLE\_64** flags are set on the **d\_map\_init** call.

```
void d_unmap_list (*handle,
*bus_list)
struct d_handle *handle;
struct dio_64 *bus_list;
```

### **Parameters**

Item	Description
<i>handle</i>	Indicates the unique handle returned by the <b>d_map_init</b> kernel service.
<i>bus_list</i>	Specifies a list of bus addresses and lengths.

### **Description**

The **d\_unmap\_list** kernel service is a bus-specific utility routine determined by the **d\_map\_init** kernel service that deallocates resources previously allocated on a **d\_map\_list** call.

The **d\_unmap\_list** kernel service must be called after I/O completion involving the area mapped by the prior **d\_map\_list** call. Some device drivers might choose to leave pages mapped for a long-term mapping of certain memory buffers. In this case, the driver must call **d\_unmap\_list** when it no longer needs the long-term mapping.

**Note:** You can use the **D\_UNMAP\_LIST** macro provided in the **/usr/include/sys/dma.h** file to code calls to the **d\_unmap\_list** kernel service. If not, you must ensure that the **d\_unmap\_list** function pointer is non-NULL before attempting the call. Not all platforms require the unmapping service.

#### **Related reference:**

“**d\_map\_init** Kernel Service” on page 82

“**d\_map\_list** Kernel Service” on page 84

## **d\_unmap\_slave Kernel Service**

### **Purpose**

Deallocates resources previously allocated on a **d\_map\_slave** call.

### **Syntax**

```
#include <sys/dma.h>
```

```
int d_unmap_slave (*handle)
struct d_handle *handle;
```

### **Parameters**

Item	Description
<i>handle</i>	Indicates the unique handle returned by the <code>d_map_init</code> kernel service.

## Description

The `d_unmap_slave` kernel service deallocates resources previously allocated on a `d_map_slave` call, disables the physical DMA channel, and returns error and status information following the DMA transfer. The `d_unmap_slave` kernel service is not an exported kernel service, but a bus-specific utility routine that is determined by the `d_map_init` kernel service and provided to the caller through the `d_handle` structure.

**Note:** You can use the `D_UNMAP_SLAVE` macro provided in the `/usr/include/sys/dma.h` file to code calls to the `d_unmap_slave` kernel service. If not, you must ensure that the `d_unmap_slave` function pointer is non-NULL before attempting to call. Not all platforms require the unmapping service.

The device driver must call `d_unmap_slave` when the I/O is complete involving a prior mapping by the `d_map_slave` kernel service.

**Note:** The `d_unmap_slave` kernel should be paired with a previous `d_map_slave` call. Multiple outstanding slave DMA transfers are not supported. This kernel service assumes that there is no DMA in progress on the affected channel and deallocates the current channel mapping.

## Return Values

Item	Description
DMA_SUCC	Indicates successful transfer. The DMA controller did not report any errors and that the Terminal Count was reached.
DMA_TC_NOTREACHED	Indicates a successful partial transfer. The DMA controller reported the Terminal Count reached for the intended transfer as set up by the <code>d_map_slave</code> call. Block devices consider this an error; however, for variable length devices this may not be an error.
DMA_FAIL	Indicates that the transfer failed. The DMA controller reported an error. The device driver assumes the transfer was unsuccessful.

## Related reference:

“`d_map_init` Kernel Service” on page 82

## `d_unmap_page` Kernel Service Purpose

Deallocates resources previously allocated on a `d_unmap_page` call.

## Syntax

```
#include <sys/dma.h>
```

```
void d_unmap_page (*handle, *busaddr)
struct d_handle *handle
uint *busaddr
```

**Note:** The following is the interface definition for `d_unmap_page` when the `DMA_ADDRESS_64` and `DMA_ENABLE_64` flags are set on the `d_map_init` call.

```
int d_unmap_page(*handle,
*busaddr)
struct d_handle *handle;
unsigned long long *busaddr;
```

## Parameters

Item	Description
<i>handle</i>	Indicates the unique handle returned by the <b>d_map_init</b> kernel service.
<i>busaddr</i>	Points to the <i>busaddr</i> field.

## Description

The **d\_unmap\_page** kernel service is a bus-specific utility routine determined by the **d\_map\_init** kernel service that deallocates resources previously allocated on a **d\_map\_page** call for a DMA master device.

The **d\_unmap\_page** service must be called after I/O completion involving the area mapped by the prior **d\_map\_page** call. Some device drivers might choose to leave pages mapped for a long-term mapping of certain memory buffers. In this case, the driver must call **d\_unmap\_page** when it no longer needs the long-term mapping.

**Note:** You can use the **D\_UNMAP\_PAGE** macro provided in the `/usr/include/sys/dma.h` file to code calls to the **d\_unmap\_page** kernel service. If not, you must ensure that the **d\_unmap\_page** function pointer is non-NULL before attempting the call. Not all platforms require the unmapping service.

### Related reference:

“d\_map\_init Kernel Service” on page 82

## dr\_reconfig System Call

### Purpose

Determines the nature of the DLPAR request.

### Syntax

```
#include <sys/dr.h>
```

```
int dr_reconfig (flags, dr_info)
int flags;
dr_info_t *dr_info;
```

### Description

The **dr\_reconfig** system call is used by DLPAR-aware applications to adjust their use of resources in relation to a DLPAR request. Applications are notified about the usage through the **SIGRECONFIG** signal, which is generated three times for each DLPAR event.

The first time is to check with the application whether the DLPAR event should be continued. Using the **DR\_EVENT\_FAIL** flag, an application can indicate that the operation should be aborted, if it is not DLPAR-safe and its operation is considered vital to the system.

The application is notified the second time before the resource is added or removed, and the third time afterwards. Applications must attempt to control their scheduling priority and policy to guarantee timely delivery of signals. The system does not guarantee that every signal that is sent is delivered before advancing to the next step in the algorithm.

The **dr\_reconfig** system call can also be used to notify applications about the changes to the workload partition that they are running. Applications are notified about changes to the CPU, memory capacity, and resources set.

The **dr\_reconfig** interface is signal-handler safe and can be used by multi-threaded programs.

The **dr\_info** structure is declared within the address space of the application. The kernel fills out data in this structure relative to the current DLPAR request. The user passes this structure identifying the current

DLPAR request, as a parameter to the kernel when the `DR_RECONFIG_DONE` flag is used. The `DR_RECONFIG_DONE` flag is used by the application to notify the kernel that necessary action to adjust their use of resources has been taken in response to the `SIGRECONFIG` signal sent to them. It is expected that the signal handler associated with the `SIGRECONFIG` signal calls the interface with the `DR_QUERY` flag to identify the phase of the DLPAR event, takes the appropriate action, and calls the interface with the `DR_RECONFIG_DONE` flag to indicate to the kernel that the signal has been handled. This type of acknowledgment to the kernel in each of the DLPAR phases enables a DLPAR event to perform efficiently.

With the addition of new fields to the `dr_info` structure, DR-aware applications can support the Micro-Partitioning<sup>®</sup> feature.

The `bindproc`, `softpset`, and `hardpset` bits are only set, if the request is to remove a cpu. If the `bindproc` is set, the process or one of its threads has a **bindprocessor** attachment, which must be resolved. If the `softpset` bit is set, the process has a Workload Manager (WLM) attachment, which can be changed by calling the appropriate WLM interface or by invoking the appropriate WLM command. If the `hardpset` bit is set, the appropriate `pset` API must be used.

**Note:** The `bcpu` and `lcpu` fields identify the cpu being removed and do not necessarily indicate that the process has a dependency that must be resolved. The `bindproc`, `softpset`, and `hardpset` bits are provided for that purpose.

The `plock` and `pshm` bits are only set, if the request is to remove memory and the process has **plock** memory or is attached to a pinned shared memory segment. If the `plock` bit is set, the process calls **plock** to unpin itself. If the `pshm` bit is set, the application has pinned shared memory segments, which may need to be detached. The memory remove request can succeed in any case, if there is enough pinnable memory in the system, so an action in this case is not necessarily required. The field `sys_pinnable_frames` provides this information, however, this value and other statistical values are just approximations. They reflect the state of the system at the time of the request. They are not updated during the request. The current size of physical memory can be determined by referencing the `_system_configuration.physmem` field.

To provide support for virtual real memory related DR operations, a new field, `dr_op`, has been added to the `dr_info` structure. The `dr_op` field provides information about the current DR operation. Additionally, all future DR operations use this field and the previously used resource bits will no longer be extended.

### **dr\_wlm\_info Structure**

```
typedef struct dr_wlm_info {
    unsigned int cpu_add : 1; // cpu wlm resource add for the WPAR
    unsigned int cpu_rem : 1; // cpu wlm resource remove for the WPAR
    unsigned int mem_add : 1; // memory wlm resource add for the WPAR
    unsigned int mem_rem : 1; // memory wlm resource remove for the WPAR
    unsigned int rs_cpu : 1; // wlm cpu rset change for the WPAR
    unsigned int rs_mem : 1; // wlm memory rset change for the WPAR
    unsigned int pad1 : 2; // un-used
    unsigned int cpu_cap : 8; // percentage of cpu capacity of the WPAR
    unsigned int mem_cap : 8; // percentage of the memory capacity of the WPAR
    unsigned int pad2 : 8; // un-used
} dr_wlm_info_t;
```

### **dr\_info Structure**

```
typedef struct dr_info {
    unsigned int add : 1; // add request
    unsigned int rem : 1; // remove request
    unsigned int cpu : 1; // target resource is a cpu
    unsigned int mem : 1; // target resource is memory
    unsigned int check : 1; // check phase in effect
    unsigned int pre : 1; // pre phase in effect
    unsigned int post : 1; // post phase in effect
    unsigned int posterror : 1; // post error phase in effect
}
```

```

    force : 1; // force option is in effect
    bindproc : 1; // process has bindprocessor dependency
    softpset : 1; // process has WLM software partition dependency
    hardpset : 1; // process has processor set API dependency
    plock : 1; // process has plock'd memory
    pshm : 1; // process has pinned shared memory
    ent_cap : 1; // target resource:entitled capacity
    var_wgt : 1; // target resource:variable weight
    splpar_capable : 1; // 1/0 partition is/not splpar capable
    splpar_shared : 1; // 1/0 partition shared/dedicated mode
    splpar_capped : 1; // 1/0 partition capped/uncapped mode
    splpar_constrained : 1; // Set to 1 if requested capacity
        update is constrained by PHYP to
        be within partition capacity bounds.
        //

unsigned int migrate : 1; // migration operation
unsigned int hibernate : 1; // hibernation operation
unsigned int partition : 1; // resource is partition
unsigned int topology_update : 1; // topology update

// The following fields are filled out for cpu based requests
int lcpu; // logical cpu ID being added or removed
int bcpu; // bind cpu ID being added or removed

// The following fields are filled out for memory based requests
size64_t req_memsz_change; // User request size in bytes
size64_t sys_memsz; // System Memory size at time of request
rpn64_t sys_free_frames; // Number of free frames in the system
rpn64_t sys_pinnable_frames; // Number of pinnable frames in system
rpn64_t sys_total_frames; // Total number of frames in system

// SPLPAR parameters.
uint64_t capacity; // partition current entitled capacity
                if ent_cap bit is set, partition's
                current variable capacity weight
                if var_wgt bit is set.
                //

int delta_cap; // delta capacity added/removed to
                current value depending on add/rem
                bit flag value above
                //

dr_wlm_info_t dr_wlm; // DR info for the WPAR
ushort dr_op; // type of DR operation

ushort dr_pad; // reserved pad field
size64_t mem_capacity; // partition's entitled
I/O memory or variable capacity.
ssize64_t delta_mem_capacity; // amount of I/O being added/removed

int reserved[2];
} dr_info_t;

```

## Parameters

**Item**  
*flags*

**Description**

The following values are supported:

**DR\_QUERY**

Identifies the current DLPAR request and the actions that the application must take to comply with the current DLPAR request. This information is returned to the caller in the structure identified by the *dr\_info* parameter.

**DR\_EVENT\_FAIL**

Fail the current DLPAR event. Root authority is required.

**DR\_RECONFIG\_DONE**

This flag is used with the **DR\_QUERY** flag. The application notifies the kernel that the actions it took to comply with the current DLPAR request are now complete. The **dr\_info** structure identifying the DLPAR request that was returned is passed as an input parameter.

*dr\_info*

Contains the address of a **dr\_info** structure, which is declared with the address space of the application.

## Return Values

Upon success, the **dr\_reconfig** system call returns a zero. If unsuccessful, it returns negative one and sets the **errno** variable to the appropriate error value.

## Error Codes

**Item**

EINVAL

ENXIO

EPERM

EINPROGRESS

**Description**

Invalid flags.

No DLPAR event in progress.

Root authority required for DR\_EVENT\_FAIL.

Cancellation of DLPAR event may only occur in the check phase.

## Related information:

Making Programs DLPAR-Aware Using DLPAR APIs

## e

The following kernel services begin with the with the letter e.

## e\_assert\_wait Kernel Service

### Purpose

Asserts that the calling kernel thread is going to sleep.

### Syntax

```
#include <sys/sleep.h>
```

```
void e_assert_wait ( event_word, interruptible)  
tid_t *event_word;  
boolean_t interruptible;
```

### Parameters



<b>Item</b>	<b>Description</b>
<i>event_word</i>	Specifies the shared event word. The kernel uses the <i>event_word</i> parameter as the anchor to the list of threads waiting on this shared event.
<i>interruptible</i>	Specifies if the sleep is interruptible.

## Description

The **e\_assert\_wait** kernel service asserts that the calling kernel thread is about to be placed on the event list anchored by the *event\_word* parameter. The *interruptible* parameter indicates whether the sleep can be interrupted.

This kernel service gives the caller the opportunity to release multiple locks and sleep atomically without losing the event should it occur. This call is typically followed by a call to either the **e\_clear\_wait** or **e\_block\_thread** kernel service. If only a single lock needs to be released, then the **e\_sleep\_thread** kernel service should be used instead.

The **e\_assert\_wait** kernel service has no return values.

## Execution Environment

The **e\_assert\_wait** kernel service can be called from the process environment only.

### Related reference:

“e\_clear\_wait Kernel Service” on page 114

“e\_sleep\_thread Kernel Service” on page 117

### Related information:

Process and Exception Management Kernel Services

## e\_block\_thread Kernel Service

### Purpose

Blocks the calling kernel thread.

### Syntax

```
#include <sys/sleep.h>
int e_block_thread ()
```

### Description

The **e\_block\_thread** kernel service blocks the calling kernel thread. The thread must have issued a request to sleep (by calling the **e\_assert\_wait** kernel service). If it has been removed from its event list, it remains runnable.

## Execution Environment

The **e\_block\_thread** kernel service can be called from the process environment only.

## Return Values

The **e\_block\_thread** kernel service return a value that indicate how the thread was awakened. The following values are defined:

Item	Description
THREAD_AWAKENED	Denotes a normal wakeup; the event occurred.
THREAD_INTERRUPTED	Denotes an interruption by a signal.
THREAD_TIMED_OUT	Denotes a timeout expiration.
THREAD_OTHER	Delineates the predefined system codes from those that need to be defined at the subsystem level. Subsystem should define their own values greater than or equal to this value.

#### Related reference:

“e\_assert\_wait Kernel Service” on page 112

#### Related information:

Process and Exception Management Kernel Services

## e\_clear\_wait Kernel Service

### Purpose

Clears the wait condition for a kernel thread.

### Syntax

```
#include <sys/sleep.h>
```

```
void e_clear_wait ( tid, result)
tid_t tid;
int result;
```

### Parameters

Item	Description
<i>tid</i>	Specifies the kernel thread to be awakened.
<i>result</i>	Specifies the value returned to the awakened kernel thread. The following values can be used:
THREAD_AWAKENED	Usually generated by the <b>e_wakeup</b> or <b>e_wakeup_one</b> kernel service to indicate a normal wakeup.
THREAD_INTERRUPTED	Indicates an interrupted sleep. This value is usually generated by a signal delivery when the <b>INTERRUPTIBLE</b> flag is set.
THREAD_TIMED_OUT	Indicates a timeout expiration.
THREAD_OTHER	Delineates the predefined system codes from those that need to be defined at the subsystem level. Subsystem should define their own values greater than or equal to this value.

### Description

The **e\_clear\_wait** kernel service clears the wait condition for the kernel thread specified by the *tid* parameter, and the thread is made runnable.

This kernel service differs from the **e\_wakeup**, **e\_wakeup\_one**, and **e\_wakeup\_w\_result** kernel services in the fact that it assumes the identity of the thread to be awakened. This kernel service should be used to handle exceptional cases, where a special action needs to be taken. The *result* parameter is used to specify the value returned to the awakened thread by the **e\_block\_thread** or **e\_sleep\_thread** kernel service.

The **e\_clear\_wait** kernel service has no return values.

## Execution Environment

The `e_clear_wait` kernel service can be called from either the process environment or the interrupt environment.

### Related reference:

“`e_wakeup`, `e_wakeup_one`, or `e_wakeup_w_result` Kernel Service” on page 121

“`e_block_thread` Kernel Service” on page 113

### Related information:

Process and Exception Management Kernel Services

## `e_sleep` Kernel Service

### Purpose

Forces the calling kernel thread to wait for the occurrence of a shared event.

### Syntax

```
#include <sys/types.h> #include <sys/errno.h> #include <sys/sleep.h> int e_sleep ( event_word, flags)
tid_t *event_word; int flags;
```

### Parameters

Item	Description
<i>event_word</i>	Specifies the shared event word. The kernel uses the <i>event_word</i> parameter to anchor the list of processes sleeping on this event. The <i>event_word</i> parameter must be initialized to <code>EVENT_NULL</code> before its first use.
<i>flags</i>	Specifies the flags that control action on occurrence of signals. These flags can be found in the <code>/usr/include/sys/sleep.h</code> file. The <i>flags</i> parameter is used to control how signals affect waiting for an event. The following flags are available to the <code>e_sleep</code> service:  <b>EVENT_SIGRET</b> Indicates the termination of the wait for the event by an unmasked signal. The return value is set to <code>EVENT_SIG</code> .  <b>EVENT_SIGWAKE</b> Indicates the termination of the event by an unmasked signal. This flag results in the transfer of control to the return from the last <code>setjmpx</code> service with the return value set to <code>EINTR</code> .  <b>EVENT_SHORT</b> Prohibits the wait from being terminated by a signal. This flag should only be used for short, guaranteed-to-wakeup sleeps.

### Description

The `e_sleep` kernel service is used to wait for the specified shared event to occur. The kernel places the current kernel thread on the list anchored by the *event\_word* parameter. This list is used by the `e_wakeup` service to wake up all threads waiting for the event to occur.

The anchor for the event list, the *event\_word* parameter, must be initialized to `EVENT_NULL` before its first use. Kernel extensions must not alter this anchor while it is in use.

The `e_wakeup` service does not wake up a thread that is not currently sleeping in the `e_sleep` function. That is, if an `e_wakeup` operation for an event is issued before the process calls the `e_sleep` service for the event, the thread still sleeps, waiting on the next `e_wakeup` service for the event. This implies that routines using this capability must ensure that no timing window exists in which events could be missed due to the `e_wakeup` service being called before the `e_sleep` operation for the event has been called.

**Note:** The `e_sleep` service can be called with interrupts disabled only if the event or lock word is pinned.

## Execution Environment

The `e_sleep` kernel service can be called from the process environment only.

## Return Values

Item	Description
<code>EVENT_SUCC</code>	Indicates a successful operation.
<code>EVENT_SIG</code>	Indicates that the <code>EVENT_SIGRET</code> flag is set and the wait is terminated by a signal.

### Related reference:

“`e_sleep1` Kernel Service”

“`e_wakeup`, `e_wakeup_one`, or `e_wakeup_w_result` Kernel Service” on page 121

### Related information:

Process and Exception Management Kernel Services

## `e_sleep1` Kernel Service

### Purpose

Forces the calling kernel thread to wait for the occurrence of a shared event.

### Syntax

```
#include <sys/types.h> #include <sys/errno.h> #include <sys/sleep.h> int e_sleep1 ( lock_word,
event_word, flags) int *lock_word; tid_t *event_word; int flags;
```

### Parameters

Item	Description
<i>lock_word</i>	Specifies the lock word for a conventional process lock.
<i>event_word</i>	Specifies the shared event word. The kernel uses this word to anchor the list of kernel threads sleeping on this event. This event word must be initialized to <code>EVENT_NULL</code> before its first use.
<i>flags</i>	Specifies the flags that control action on occurrence of a signal. These flags are found in the <code>/usr/include/sys/sleep.h</code> file.

### Description

**Note:** The `e_sleep1` kernel service is provided for porting old applications written for previous versions of the operating system. Use the `e_sleep_thread` kernel service when writing new applications.

The `e_sleep1` kernel service waits for the specified shared event to occur. The kernel places the current kernel thread on the list anchored by the *event\_word* parameter. The `e_wakeup` service wakes up all threads on the list.

The `e_wakeup` service does not wake up a thread that is not currently sleeping in the `e_sleep1` function. That is, if an `e_wakeup` operation for an event is issued before the thread calls the `e_sleep1` service for the event, the thread still sleeps, waiting on the next `e_wakeup` operation for the event. This implies that routines using this capability must ensure that no timing window exists in which events could be missed due to the `e_wakeup` service being called before the `e_sleep1` service for the event has been called.

The `e_sleep1` service also unlocks the conventional lock specified by the *lock\_word* parameter before putting the thread to sleep. It also reacquires the lock when the thread wakes up.

The anchor for the event list, specified by the *event\_word* parameter, must be initialized to `EVENT_NULL` before its first use. Kernel extensions must not alter this anchor while it is in use.

**Note:** The `e_sleep1` service can be called with interrupts disabled, only if the event or lock word is pinned.

### Values for the flags Parameter

The *flags* parameter controls how signals affect waiting for an event. There are three flags available to the `e_sleep1` service:

Item	Description
EVENT_SIGRET	Indicates the termination of the wait for the event by an unmasked signal. The return value is set to EVENT_SIG.
EVENT_SIGWAKE	Indicates the termination of the event by an unmasked signal. This flag also indicates the transfer of control to the return from the last <code>setjmpx</code> service with the return value set to <code>EINTR</code> .
EVENT_SHORT	Indicates that signals cannot terminate the wait. Use the <code>EVENT_SHORT</code> flag for only short, guaranteed-to-wakeup sleeps.

**Note:** The `EVENT_SIGRET` flag overrides the `EVENT_SIGWAKE` flag.

### Execution Environment

The `e_sleep1` kernel service can be called from the process environment only.

### Return Values

Item	Description
EVENT_SUCC	Indicates successful completion.
EVENT_SIG	Indicates that the <code>EVENT_SIGRET</code> flag is set and the wait is terminated by a signal.

### Related reference:

“`e_sleep` Kernel Service” on page 115

“`e_wakeup`, `e_wakeup_one`, or `e_wakeup_w_result` Kernel Service” on page 121

### Related information:

Interrupt Environment

## `e_sleep_thread` Kernel Service

### Purpose

Forces the calling kernel thread to wait for the occurrence of a shared event.

### Syntax

```
#include <sys/sleep.h>
```

```
int e_sleep_thread ( event_word, lock_word, flags)
tid_t *event_word;
void *lock_word;
int flags;
```

### Parameters

Item	Description
<i>event_word</i>	Specifies the shared event word. The kernel uses the <i>event_word</i> parameter as the anchor to the list of threads waiting on this shared event.
<i>lock_word</i>	Specifies simple or complex lock to unlock.
<i>flags</i>	Specifies lock and signal handling options.

## Description

The **e\_sleep\_thread** kernel service forces the calling thread to wait until a shared event occurs. The kernel places the calling thread on the event list anchored by the *event\_word* parameter. This list is used by the **e\_wakeup**, **e\_wakeup\_one**, and **e\_wakeup\_w\_result** kernel services to wakeup some or all threads waiting for the event to occur.

A lock can be specified; it will be unlocked when the kernel service is entered, just before the thread blocks. This lock can be a simple or a complex lock, as specified by the *flags* parameter. When the kernel service exits, the lock is re-acquired.

## Flags

The *flags* parameter specifies options for the kernel service. Several flags can be combined with the bitwise OR operator. They are described below.

The four following flags specify the lock type. If the *lock\_word* parameter is not **NULL**, exactly one of these flags must be used.

Flag	Description
<b>LOCK_HANDLER</b>	<i>lock_word</i> specifies a simple lock protecting a thread-interrupt or interrupt-interrupt critical section.
<b>LOCK_SIMPLE</b>	<i>lock_word</i> specifies a simple lock protecting a thread-thread critical section.
<b>LOCK_READ</b>	<i>lock_word</i> specifies a complex lock in shared-read mode.
<b>LOCK_WRITE</b>	<i>lock_word</i> specifies a complex lock in exclusive write mode.

The following flag specify the signal handling. By default, while the thread sleeps, signals are held pending until it wakes up.

Item	Description
<b>INTERRUPTIBLE</b>	The signals must be checked while the kernel thread is sleeping. If a signal needs to be delivered, the thread is awakened.

## Return Values

The **e\_sleep\_thread** kernel service return a value that indicate how the kernel thread was awakened. The following values are defined:

Item	Description
<b>THREAD_AWAKENED</b>	Denotes a normal wakeup; the event occurred.
<b>THREAD_INTERRUPTED</b>	Denotes an interruption by a signal. This value can be returned even if the <b>INTERRUPTIBLE</b> flag is not set since it may be also generated by the <b>e_clear_wait</b> or <b>e_wakeup_w_result</b> kernel services.
<b>THREAD_TIMED_OUT</b>	Denotes a timeout expiration. The <b>e_sleep_thread</b> has no timeout. However, the <b>e_clear_wait</b> or <b>e_wakeup_w_result</b> kernel services may generate this return value.
<b>THREAD_OTHER</b>	Delineates the predefined system codes from those that need to be defined at the subsystem level. Subsystem should define their own values greater than or equal to this value.

## Execution Environment

The `e_sleep_thread` kernel service can be called from the process environment only.

### Related reference:

“`e_wakeup`, `e_wakeup_one`, or `e_wakeup_w_result` Kernel Service” on page 121

“`e_block_thread` Kernel Service” on page 113

### Related information:

Locking Kernel Services

## et\_post Kernel Service

### Purpose

Notifies a kernel thread of the occurrence of one or more events.

### Syntax

```
#include <sys/types.h> #include <sys/errno.h> #include <sys/sleep.h> void et_post ( events, tid)
unsigned long events; tid_t tid;
```

### Parameters

Item	Description
<i>events</i>	Identifies the masks of events to be posted.
<i>tid</i>	Specifies the thread identifier of the kernel thread to be notified.

### Description

The `et_post` kernel service is used to notify a kernel thread that one or more events occurred.

The `et_post` service provides the fastest method of interprocess communication, although only the event numbers are passed.

The event numbers must be known by the cooperating components, either through programming convention or the passing of initialization parameters.

The `et_post` service is performed automatically when sending a request to a device queue serviced by a kernel thread or when sending an acknowledgment.

The `EVENT_KERNEL` mask defines the event bits reserved for use by the kernel. For example, a bit with a value of 1 indicates an event bit reserved for the kernel. Kernel extensions should assign their events starting with the most significant bits and working down. If threads using the `et_post` service are also using the device queue management kernel extensions, care must be taken not to use the event bits registered for device queue management.

The `et_wait` service does not sleep but returns immediately if a specified event has already been posted by the `et_post` service.

## Execution Environment

The `et_post` kernel service can be called from either the process or interrupt environment.

## Return Values

The `et_post` service has no return values.

### Related reference:

“`et_wait` Kernel Service”

### Related information:

Process and Exception Management Kernel Services

## `et_wait` Kernel Service

### Purpose

Forces the calling kernel thread to wait for the occurrence of an event.

### Syntax

```
#include <sys/types.h> #include <sys/errno.h> #include <sys/sleep.h> unsigned long et_wait (  
wait_mask, clear_mask, flags) unsigned long wait_mask; unsigned long clear_mask; int flags;
```

### Parameters

Item	Description
<i>wait_mask</i>	Specifies the mask of events to await.
<i>clear_mask</i>	Specifies the mask of events to clear.
<i>flags</i>	Specifies the flags controlling actions on occurrence of a signal.

The *flags* parameter is used to control how signals affect waiting for an event. There are two flag values:

**EVENT\_SIGRET**  
Causes the wait for the event to be ended by an unmasked signal and the return value set to **EVENT\_SIG**.

**EVENT\_SIGWAKE**  
Causes the event to be ended by an unmasked signal and control transferred to the return from the last `setjmpx` call, with the return value set to **EXSIG**.

**EVENT\_SHORT**  
Prohibits the wait from being terminated by a signal. This flag should only be used for short, guaranteed-to-wakeup sleeps.

**Note:** The **EVENT\_SIGRET** flag overrides the **EVENT\_SIGWAKE** flag.

### Description

The `et_wait` kernel service forces the calling kernel thread to wait for specified events to occur.

The *wait\_mask* parameter indicates a mask, where each bit set equal to 1 represents an event for which the thread must wait. The *clear\_mask* parameter indicates a mask of events that must clear when the wait is complete. Subsequent calls to the `et_wait` service return immediately unless you clear the bits, which ends the wait.

**Note:** The `et_wait` service can be called with interrupts disabled only if the event or lock word is pinned.

### Strategies for Using `et_wait`

Calling the `et_wait` kernel service with the **EVENT\_SIGRET** flag clears the the pending events field when the signal is received. If `et_wait` is called again by the same kernel thread, the thread waits indefinitely for an event that has already occurred. When this happens, the thread does not run to completion. This problem occurs only if the event and signal are posted at the same time.



To avoid this problem, use one of the following programming methods:

- Use the **EVENT\_SHORT** flag to prevent signals from waking the thread up.
- Mask signals prior to the call of **et\_wait** by using the **limit\_sigs** kernel service. Then call **et\_wait**. Invoke the **sigprocmask** call to restore the signal mask by using the mask returned previously by **limit\_sigs**.

The **et\_wait** service is also used to clear events without waiting for them to occur. This is accomplished by doing one of the following:

- Set the *wait\_mask* parameter to **EVENT\_NDELAY**.
- Set the bits in the *clear\_mask* parameter that correspond with the events to be cleared to 1.

Because the **et\_wait** service returns an event mask indicating those events that were actually cleared, these methods can be used to poll the events.

## Execution Environment

The **et\_wait** kernel service can be called from the process environment only.

## Return Values

Upon successful completion, the **et\_wait** service returns an event mask indicating the events that terminated the wait. If an **EVENT\_NDELAY** value is specified, the returned event mask indicates the pending events that were cleared by this call. Otherwise, it returns the following error code:

Item	Description
<b>EVENT_SIG</b>	Indicates that the <b>EVENT_SIGRET</b> flag is set and the wait is terminated by a signal.

### Related reference:

“**et\_post** Kernel Service” on page 119

“**setjmpx** Kernel Service” on page 468

### Related information:

Process and Exception Management Kernel Services

## **e\_wakeup, e\_wakeup\_one, or e\_wakeup\_w\_result** Kernel Service Purpose

Notifies kernel threads waiting on a shared event of the event's occurrence.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/sleep.h>
```

```
void e_wakeup ( event_word)
tid_t *event_word;
```

```
void e_wakeup_one ( event_word)
tid_t *event_word;
```

```
void e_wakeup_w_result ( event_word, result)
tid_t *event_word;
int result;
```

## Parameters

Item	Description
<i>event_word</i>	Specifies the shared event designator. The kernel uses the <i>event_word</i> parameter as the anchor to the list of threads waiting on this shared event.
<i>result</i>	Specifies the value returned to the awakened kernel thread. The following values can be used: <b>THREAD_AWAKENED</b> Indicates a normal wakeup. This is the value automatically generated by the <b>e_wakeup</b> or <b>e_wakeup_one</b> kernel services. <b>THREAD_INTERRUPTED</b> Indicates an interrupted sleep. This value is usually generated by a signal delivery when the <b>INTERRUPTIBLE</b> flag is set. <b>THREAD_TIMED_OUT</b> Indicates a timeout expiration. <b>THREAD_OTHER</b> Delineates the predefined system codes from those that need to be defined at the subsystem level. Subsystem should define their own values greater than or equal to this value.

## Description

The **e\_wakeup** and **e\_wakeup\_w\_result** kernel services wake up all kernel threads sleeping on the event list anchored by the *event\_word* parameter. The **e\_wakeup\_one** kernel service wakes up only the most favored thread sleeping on the event list anchored by the *event\_word* parameter.

When threads are awakened, they return from a call to either the **e\_block\_thread** or **e\_sleep\_thread** kernel service. The return value depends on the kernel service called to wake up the threads (the wake-up kernel service):

- **THREAD\_AWAKENED** is returned if the **e\_wakeup** or **e\_wakeup\_one** kernel service is called
- The value of the *result* parameter is returned if the **e\_wakeup\_w\_result** kernel service is called.

If a signal is delivered to a thread being awakened by one of the wake-up kernel services, and if the thread specified the **INTERRUPTIBLE** flag, the signal delivery takes precedence. The thread is awakened with a return value of **THREAD\_INTERRUPTED**, regardless of the called wake-up kernel service.

The **e\_wakeup** and **e\_wakeup\_w\_result** kernel services set the *event\_word* parameter to **EVENT\_NULL**.

The **e\_wakeup**, **e\_wakeup\_one**, and **e\_wakeup\_w\_result** kernel services have no return values.

## Execution Environment

The **e\_wakeup**, **e\_wakeup\_one**, and **e\_wakeup\_w\_result** kernel services can be called from either the process environment or the interrupt environment.

When called by an interrupt handler, the *event\_word* parameter must be located in pinned memory.

### Related reference:

“e\_clear\_wait Kernel Service” on page 114

“e\_sleep\_thread Kernel Service” on page 117

### Related information:

Process and Exception Management Kernel Services

## **e\_wakeup\_w\_sig** Kernel Service Purpose

Posts a signal to sleeping kernel threads.

## Syntax

```
#include <sys/sleep.h>
```

```
void e_wakeup_w_sig ( event_word, sig)
tid_t *event_word;
int sig;
```

## Parameters

Item	Description
<i>event_word</i>	Specifies the shared event word. The kernel uses the <i>event_word</i> parameter as the anchor to the list of threads waiting on this shared event.
<i>sig</i>	Specifies the signal number to post.

## Description

The **e\_wakeup\_w\_sig** kernel service posts the signal *sig* to each kernel thread sleeping interruptible on the event list anchored by the *event\_word* parameter.

The **e\_wakeup\_w\_sig** kernel service has no return values.

## Execution Environment

The **e\_wakeup\_w\_sig** kernel service can be called from either the process environment or the interrupt environment.

### Related reference:

“e\_block\_thread Kernel Service” on page 113

“e\_clear\_wait Kernel Service” on page 114

### Related information:

Process and Exception Management Kernel Services

## eeh\_broadcast Kernel Service

### Purpose

This service is provided for device drivers to coordinate activities during an EEH event.

## Syntax

```
void eeh_broadcast(handle, message)
eeh_handle_t handle;
unsigned long long message;
```

## Parameters

Item	Description
<i>handle</i>	EEH handle obtained from <b>eeh_init</b> or <b>eeh_init_multifunc</b>
<i>message</i>	User- or kernel-defined message

## Description

Because single-function drivers do not have a need for coordination, this service is intended for multifunction drivers only. If a single-function driver calls it, it is a NOP. There are two kinds of messages that can be sent among the drivers: kernel-defined messages (such as **EEH\_DD\_SUSPEND** and **EEH\_DD\_DEAD**) and the user-defined messages. See **sys/eeh.h** for help on how to define user messages.

Kernel messages have a higher priority than user messages. Therefore, if user messages and kernel messages are both pending, the kernel messages are sent out before the user messages.

**Note:** Device drivers should only broadcast their own messages (that is, the user-defined message) and not the kernel messages.

Within the kernel messages, EEH\_DD\_DEAD has the highest priority. Multiple messages of the same kind may or may not be coalesced depending upon the relative timing. Messages are sent by invoking the callback routines. The callback routines are invoked sequentially but not in any specific order except that the last driver to receive a message will have the EEH\_MASTER flag set to indicate that all other drivers have finished processing the message. Only one message is broadcast at a time—that is, all registered callback routines are called sequentially with the same message before moving on to the next message. Finally, they are invoked asynchronously at INTIODONE priority. Because they are broadcast asynchronously, a device driver must not assume on a specific timeout within which the message would arrive.

The macro `EEH_BROADCAST(handle, message)` is provided for device drivers to call this service.

## Execution Environment

This kernel service can be called from the process or interrupt environment.

## Return Values

This service has no return value.

### Related reference:

“`eeh_enable_slot` Kernel Service” on page 128

“`eeh_init_multifunc` Kernel Service” on page 130

“`eeh_slot_error` Kernel Service” on page 136

## eeh\_clear Kernel Service Purpose

This service unregisters a slot for an EEH function and removes resources allocated by the `eeh_init` or `eeh_init_multifunc` kernel service.

## Syntax

```
#include <sys/eeh.h>
```

```
void eeh_clear(handle)  
eeh_handle_t handle;
```

## Parameters

Item	Description
<i>handle</i>	EEH handle obtained from <code>eeh_init</code> or <code>eeh_init_multifunc</code> kernel services

## Description

**Single-function Drivers:** This service disables EEH function on the slot and frees its `eeh_handle`.

**Multifunction Drivers:** For a multifunction adapter driver, this service removes the driver from a list of registered drivers under the same parent bus. This service also disables EEH function on the slot if this is the last driver to unregister and the state of the slot is NORMAL.

All device drivers are required to call `eeh_clear` before being removed from the system, so that there are no hot plug conflicts. A subsequent adapter might fail in `eeh_init_multifunc()` on the slot if the `eeh_clear` kernel service has not cleared the prior device drivers on that slot. A driver can unregister at unconfigure/unload time. The kernel checks the state of the slot when this service is called. If the slot state is neither NORMAL nor DEAD, `eeh_clear` sleeps until the state returns to one of them.

The macro `EEH_CLEAR(handle)` is provided for device drivers to call this service. This service is called by a function pointer in the EEH handle.

## Execution Environment

This kernel service can only be called from the process environment.

## Return Values

This service has no return values.

### Related reference:

“`eeh_broadcast` Kernel Service” on page 123

“`eeh_enable_slot` Kernel Service” on page 128

“`eeh_init` Kernel Service” on page 129

## `eeh_disable_slot` Kernel Service Purpose

This service disables a slot for the EEH operations.

## Syntax

```
#include <sys/eeh.h>
```

```
long eeh_disable_slot(handle)  
eeh_handle_t handle;
```

## Parameters

Item	Description
<i>handle</i>	EEH handle obtained from the <code>eeh_init</code> kernel service

## Description

This service disables EEH operation on a slot.

### CAUTION:

**CAUTION: Disabling EEH operation on a slot is highly discouraged, because it can cause system crash or worse, data corruption.**

This service can only be called by the single-function adapter drivers. If the service fails for a hardware or firmware reason, an error is logged.

Multifunction drivers call this service indirectly via `eeh_clear()`. It fails with `EEH_FAIL` if called directly by a multifunction driver.

The macro `EEH_DISABLE_SLOT(handle)` is provided for device drivers to call this service.

## Execution Environment

This kernel service can be called from the process or interrupt environment.

## Return Values

Item	Description
EEH_SUCC	Slot successfully disabled
EEH_FAIL	Unable to disable the slot

### Related reference:

“`eeh_enable_slot` Kernel Service” on page 128

“`eeh_read_slot_state` Kernel Service” on page 133

“`eeh_slot_error` Kernel Service” on page 136

## eeh\_enable\_dma Kernel Service

### Purpose

This service enables DMA operations to an adapter after an EEH event.

### Syntax

```
#include <sys/eeh.h>
```

```
long eeh_enable_dma(handle)  
eeh_handle_t handle;
```

### Parameters

Item	Description
<i>handle</i>	EEH handle obtained from <code>theeh_init</code> or <code>eeh_init_multifunc</code> kernel services

### Description

When an EEH event occurs on a slot, all Direct Memory Access (DMA) operations on the slot are inhibited. This service should be called to re-enable DMA after an EEH event. This service can only be called from the dump context (that is, when the dump is in progress).

**Single-function Drivers:** This service enables the DMA operations on a slot. If this call fails with `EEH_FAIL`, an error is logged by the kernel.

**Multifunction Drivers:** On the multifunction adapters, the slot state must be either `SUSPEND` or `DEBUG`, and the caller must be an `EEH_MASTER`. This service is called only from a dump context. While a system dump is in progress, all callbacks and broadcasts are suspended, and a multifunction adapter is treated like a single-function adapter, because the system can no longer support the EEH multifunction kernel services. If the service fails, `EEH_FAIL` is returned. If the failure is due to hardware or firmware, an error is logged.

There are cases when this kernel service cannot succeed because of the platform state restrictions. In such a case, if a driver calls it, the service would return `EEH_FAIL`. This causes the slot to be marked permanently unavailable, which is not correct because the slot can be recovered. To avoid receiving `EEH_FAIL` from this service, the driver should supply the `EEH_ENABLE_NO_SUPPORT_RC` flag at `eeh_init_multifunc()` time. If the `EEH_ENABLE_NO_SUPPORT_RC` flag is supplied, `eeh_enable_dma()` returns `EEH_NO_SUPPORT`, indicating to the drivers that they cannot collect debug data but must continue with the next step in recovery.

The macro `EEH_ENABLE_DMA(handle)` is provided for device drivers to call this service.

## Execution Environment

This kernel service can only be called from a process or interrupt environment.

## Return Values

This kernel service has no return values.

### Related reference:

“`eeh_disable_slot` Kernel Service” on page 125

“`eeh_enable_pio` Kernel Service”

“`eeh_enable_slot` Kernel Service” on page 128

## `eeh_enable_pio` Kernel Service

### Purpose

This kernel service enables programmed I/O (PIO or MMIO) to an adapter after an EEH event.

### Syntax

```
#include <sys/eeh.h>
```

```
long eeh_enable_pio(handle)
eeh_handle_t handle;
```

### Parameters

Item	Description
<i>handle</i>	EEH handle obtained from the <code>eeh_init</code> or <code>eeh_init_multifunc</code> kernel services

### Description

When an EEH event occurs on a slot, all load and store operations (such as PIO) are inhibited. This kernel service should be called to re-enable PIO after an EEH event.

**Single-function Drivers:** This kernel service enables the load and store operations on a slot. If this call fails with `EEH_FAIL`, an error is logged by the kernel.

**Multifunction Drivers:** On the multifunction adapters, the state of the slot is checked for either `SUSPEND` or `DEBUG`. The caller must be an `EEH_MASTER`. If the state is `SUSPEND`, a series of device driver callback routines is executed with a command option of `EEH_DD_DEBUG` and flag set to `EEH_DD_PIO_ENABLED`. The callbacks inform device drivers that PIO has been enabled and that further debug procedures can be executed (such as reading command and status register). This service can be called as a result of the `EEH_DD_SUSPEND` or `EEH_DD_DEBUG` callback message as many times as needed by the `EEH_MASTER`. Additional calls to this service trigger a new set of callbacks. If this service fails, `EEH_FAIL` is returned. If the failure is due to hardware or firmware, an error is logged.

There are cases when this kernel service cannot succeed due to the platform state restrictions. In such a case, if a driver calls it, the kernel service would return `EEH_FAIL` followed by a `EEH_DD_DEAD` message. This causes the slot to be marked permanently unavailable, which is not correct because the slot can be recovered. To avoid receiving `EEH_FAIL` from this service, the driver should supply the `EEH_ENABLE_NO_SUPPORT_RC` flag at `eeh_init_multifunc()` time. If the `EEH_ENABLE_NO_SUPPORT_RC` flag is supplied, `eeh_enable_pio()` returns `EEH_NO_SUPPORT`, indicating to the drivers that they cannot collect debug data but must continue with the next step in recovery.

The macro `EEH_CLEAR(handle)` is provided for device drivers to call this service. This service is called via a function pointer in the EEH handle.

**Note:** Enabling PIO is not the same as recovering the slot. In fact, this is an optional step in the recovery procedure.

## Execution Environment

This kernel service can be called from the process or interrupt environment.

## Return Values

Item	Description
EEH_SUCC	PIO successfully enabled.
EEH_FAIL	Invalid call or could not enable PIO.
EEH_NO_SUPPORT	Call is valid according to AIX EEH state, but current platform state precludes normal completion.

### Related reference:

“`eeh_disable_slot` Kernel Service” on page 125

“`eeh_enable_dma` Kernel Service” on page 126

“`eeh_enable_slot` Kernel Service”

## eeh\_enable\_slot Kernel Service

### Purpose

This service enables a slot for the EEH operations.

### Syntax

```
#include <sys/eeh.h>
```

```
long eeh_enable_slot(handle)  
eeh_handle_t handle;
```

### Parameters

Item	Description
<i>handle</i>	EEH handle obtained from <code>theeh_init</code> kernel service

### Description

This service enables EEH operation on a slot so that when certain errors occur on a PCI bus, the slot will freeze (that is, PIO and DMA are disabled, which prevents potential system crash, data corruption, and so on). This service can only be called by the single-function adapter drivers. If the service fails for hardware or firmware reasons, an error is logged.

Multifunction drivers call this service indirectly via `eeh_init_multifunc()`. It fails with `EEH_FAIL` if called directly by a multifunction driver.

The macro `EEH_ENABLE_SLOT(handle)` is provided for device drivers to call this service.

## Execution Environment

This kernel service can be called from the process or interrupt environment.

## Return Values



Item	Description
EEH_SUCC	Slot successfully enabled
EEH_FAIL	Unable to enable the slot

#### Related reference:

“`eeh_disable_slot` Kernel Service” on page 125

“`eeh_enable_dma` Kernel Service” on page 126

“`eeh_enable_pio` Kernel Service” on page 127

## eeh\_init Kernel Service

### Purpose

This service registers a single-function adapter slot on a PCI/PCI-E bus for EEH function.

### Syntax

```
#include <sys/eeh.h>
```

```
eeh_handle_t eeh_init(pbid, slot, flag)
long      pbid;
long      slot;
long      flag;
```

### Parameters

Item	Description
<i>pbid</i>	AIX parent bus identifier
<i>slot</i>	device slot (device*8+function). This is same as "connwhere" property in CuDv.
<i>flag</i>	flag that enables eeh

### Description

The *pbid* argument identifies a bus type and number. The bus type is `IO_PCI` in the case of PCI and PCI-X bus. If the bus type is `IO_PCIE`, the device is on PCI-E (PCI Express) bus. The bus number is a unique identifier determined during bus configuration. The `BID_VAL` macro defined in `ioacc.h` is used to generate the *bid*. The *slot* argument is the device/function combination ((device\*8) + function) as in the PCI addressing scheme. The *flag* argument of `EEH_ENABLE` enables the slot. The *flag* argument of `EEH_DISABLE` does not enable the slot but still allocates an EEH handle. This service should be called only by the single-function adapter drivers.

The macro `EEH_INIT(pbid, slot, flag)` is provided for the device drivers to call this service. The `eeh_handle` is defined as follows in `<sys/eeh.h>`:

```
/*
 * This is the eeh_handle structure for the eeh_* services
 */
typedef struct eeh_handle *      eeh_handle_t;
struct eeh_handle {
    struct eeh_handle *next;
    long      bid;                /* bus id passed to eeh_init */
    long      slot;              /* slot passed to eeh_init */
    long      flag;              /* flag passed to eeh_init */
    int       config_addr;       /* Configuration Space Address */
    int       eeh_mode;          /* Indicates safe mode */
    uint      retry_delay;       /* re-read the slot state after *
                                * these many seconds. */
    int       reserved1;
    int       reserved2;
    int       reserved3;
    long long PHB_Unit_ID;       /* /pci@ */
};
```

```

void    (*eeh_clear)(eeh_handle_t);
long    (*eeh_enable_pio)(eeh_handle_t);
long    (*eeh_enable_dma)(eeh_handle_t);
long    (*eeh_reset_slot)(eeh_handle_t, int);
long    (*eeh_enable_slot)(eeh_handle_t);
long    (*eeh_disable_slot)(eeh_handle_t);
long    (*eeh_read_slot_state)(eeh_handle_t, long *, long *);
long    (*eeh_slot_error)(eeh_handle_t, int, char *, long);
struct eeh_shared_domain *parent_sd; /* point back to the parent
                                     * shared domain structure if
                                     * in shared domain, NULL if singlefunc.
                                     */
void    (*eeh_configure_bridge)(eeh_handle_t);
void    (*eeh_broadcast)(eeh_handle_t, unsigned long long);
};

```

This is an exported kernel service.

## Execution Environment

This service can only be called from the process environment.

## Return Values

Item	Description
EEH_FAIL	Unable to allocate EEH handle.
EEH_NO_SUPPORT	EEH not supported on this system, no handle allocated.
struct eeh_handle *	If successful.

## Related reference:

“eeh\_broadcast Kernel Service” on page 123

“eeh\_clear Kernel Service” on page 124

“eeh\_enable\_slot Kernel Service” on page 128

## eeh\_init\_multifunc Kernel Service

### Purpose

This kernel service registers a multifunction adapter slot on a PCI/PCI-E bus for EEH function.

### Syntax

```

#include <sys/eeh.h>

eeh_handle_t eeh_init_multifunc(gpbid, pbid, slot, flag, delay_seconds,
                               callback_ptr, dds_ptr)

long gpbid;
long pbid;
long slot;
long flag;
long delay_seconds;
long (*callback_ptr)();
void *dds_ptr;

```

### Parameters

Item	Description
<i>gpbid</i>	Bus identifier of grandparent bus.
<i>pbid</i>	Bus identifier of parent bus.
<i>slot</i>	Slot on the parent bus (device*8+function). This is same as "connwhere" property in CuDv for the device.
<i>flag</i>	Flag that enables eeh, checks if the slot is already taken, etc.
<i>delay_seconds</i>	Time delay after a reset (in seconds).
<i>callback_ptr</i>	Device driver callback routine.
<i>dds_ptr</i>	Cookie to a target device driver that is usually a pointer to the adapter structure.

## Description

This kernel service is provided for systems that support shared EEH domain, where one or more PCI functions in one or more adapters could belong to the same EEH recovery domain. In the past, this was called "multifunction adapter". The shared EEH domain is a more general concept than just a multifunction adapter. It is also recommended that single function adapters use the shared EEH model. All PCI-E devices, single or multifunction have to use the shared EEH model and hence this kernel service to register for EEH (instead of `eeh_init(0)`). In a shared EEH domain, multiple instances of device drivers may be operating. The instances are independent of each other and hence oblivious to each other's existence. Therefore, when recovering a slot from an EEH event, there is a need to coordinate the recovery procedure among them. As with `eeh_init(0)`, this service also returns an `eeh_handle` to the calling device driver.

There are two kinds of adapters: bridged and non-bridged. A bridged adapter has a bridge on the card such as PCI-to-PCI or PCIX-to-PCIX or PCI-E switch. For PCI and PCI-X bridged-adapters, *pbid* is the bus ID of the parent bus, and *gpbid* is the bus ID of the grandparent bus. The parent bus for a bridged adapter is the bus generated by the bridge/switch on the adapter. A *bid* identifies a bus number and type. The bus type is `IO_PCI` in the case of PCI and PCI-X bus, and `IO_PCIE` in the case of PCI-E bus. The bus number is a unique identifier determined during bus configuration. The `BID_VAL` macro defined in `ioacc.h` is used to generate the *bid*. For non-bridged adapters, *pbid* and *gpbid* are the same and are the bus IDs of the parent bus. Thus, when *pbid* and *gpbid* have different values for a PCI or PCI-X device, the kernel knows that this is a bridged adapter and needs to the bridge recovered as part of EEH recovery. It is not necessary to know if a PCI-E device is bridged or not for the purposes of EEH. Therefore, *pbid* and *gpbid* must be same and equal to the parent bus bid.

In summary, there are the following cases:

1. PCI/PCI-X non-bridged adapters and all PCI-E adapters: *gpbid* and *pbid* are same and equal to the parent bus *bid*.
2. PCI/PCI-X bridged adapters, *gpbid* is grandparent bus *bid*, and *pbid* is parent bus bid.

The *slot* argument is the device/function combination ((device\* 8) + function) as in the PCI addressing scheme. This is the same as the `connwhere` ODM value of the device.

The following flag values are legal:

Item	Description
<code>EEH_ENABLE_FLAG/EEH_DISABLE_FLAG</code>	The slot is always enabled for EEH when this service is called by the first driver on that slot. All subsequent requests to enable the slot via the <code>EEH_ENABLE</code> flag are ignored. Therefore, the flag argument of <code>EEH_ENABLE</code> is optional, and a flag of <code>EEH_DISABLE</code> is ignored.
<code>EEH_CHECK_SLOT</code>	The flag argument of <code>EEH_CHECK_SLOT</code> verifies whether a given slot is already registered. A value of either <code>EEH_SLOT_ACTIVE</code> or <code>EEH_SLOT_FREE</code> is returned. No registration occurs with the <code>EEH_CHECK_SLOT</code> flag, and it supersedes all other flags. This flag simply checks the slot and returns without any other action.

Item	Description
EEH_ENABLE_NO_SUPPORT_RC	If the flag is set to <code>EEH_ENABLE_NO_SUPPORT_RC</code> , <code>eeh_enable_pio()</code> and <code>eeh_enable_dma()</code> return <code>EEH_NO_SUPPORT</code> under certain conditions. See “ <code>eeh_enable_dma</code> Kernel Service” on page 126 and “ <code>eeh_enable_pio</code> Kernel Service” on page 127 for more information.

Multiple flags can be logically ORed together.

The slot is always enabled for EEH when this service is called by the first driver on that slot. All subsequent requests to enable the slot via the `EEH_ENABLE` flag are ignored. Therefore, the *flag* argument of `EEH_ENABLE` is optional, and a flag of `EEH_DISABLE` is ignored. The flag argument of `EEH_CHECK_SLOT` verifies whether a given slot is already registered. A value of either `EEH_SLOT_ACTIVE` or `EEH_SLOT_FREE` is returned. No registration will occur with the `EEH_CHECK_SLOT` flag, and it supersedes all other flags. This flag just checks the slot and returns without any other action. If the flag is set to `EEH_ENABLE_NO_SUPPORT_RC`, `eeh_enable_pio()` and `eeh_enable_dma()` returns `EEH_NO_SUPPORT` under certain conditions. See `eeh_enable_pio()` and `eeh_enable_dma()` for more information. It is allowed to logically OR multiple flags together.

The *delay\_seconds* argument allows the device driver to set a time delay between completion of PCI reset and configuration of the bridge on the adapter. The delay is enforced even if the adapter is non-bridged. If a value of 0 is specified for *delay\_seconds*, a default delay time of 1 second is set. When several drivers register on the same *pbid* (under a shared EEH domain), the highest delay time among all registered drivers is used.

The *callback\_ptr* argument is a function pointer to an EEH callback routine. The handler is defined by the device driver and is called by the kernel in order to coordinate recovery among different drivers on the same slot. The driver handles a variety of messages from the kernel in its callback routine. These messages trigger the next step in recovery. The callback routines are called sequentially at INTIODONE interrupt level.

The *dds\_ptr* argument is a cookie that is passed to the driver when the callback routine is invoked. Drivers normally specify a pointer to the device driver's adapter structure.

**EEH\_SAFE mode:** A bridged adapter needs to have its bridge reconfigured at the end of PCI reset. However, if the platform firmware does not support reconfiguration of the bridge, the adapter is marked as `EEH_SAFE` by the kernel. An `EEH_SAFE` adapter cannot finish error recovery after an EEH event because of the unsatisfied firmware dependency. See `eeh_reset_slot` for information on how the error recovery is handled in `EEH_SAFE` mode.

The macro `EEH_INIT_MULTIFUNC(gpbid, pbid, slot, flag, delay_seconds, callback_ptr, dds_ptr)` is provided for the device drivers in order to call this service. This is an exported kernel service.

## Execution Environment

This kernel service can only be called from the process environment.

## Return Values

Item	Description
EEH_FAIL	Unable to allocate EEH handle.
EEH_NO_SUPPORT	EEH is not supported on this system, no handle allocated.
EEH_SLOT_ACTIVE	Given slot is already registered.
EEH_SLOT_FREE	Given slot free.
EEH_BUSY	Unable to continue, because the slot is in the middle of error recovery.
struct eeh_handle *	Upon Success.

#### Related reference:

“eeh\_broadcast Kernel Service” on page 123

“eeh\_clear Kernel Service” on page 124

## eeh\_read\_slot\_state Kernel Service

### Purpose

This service returns state and capabilities of a slot with respect to EEH operation.

### Syntax

```
long eeh_read_slot_state(handle, state, support)
eeh_handle_t handle;
long *state;
long *support;
```

### Parameters

Item	Description
<i>handle</i>	EEH handle obtained from <code>eeh_init</code> or <code>eeh_init_multifunc</code>
<i>state</i>	State of a slot with respect to EEH
<i>support</i>	Indicates if EEH is supported by this slot

### Description

This service is used to query the hardware state of a slot and to determine whether a given slot supports EEH. It should be called to confirm an EEH event if the driver suspects that the PIO data is invalid (for example, getting all Fs from reading a register). This service returns the hardware state in *state* and indicates whether the slot supports EEH in *support*. The *state* and *support* parameters are integer values as shown below:

Valid *state* values are as follows:

Item	Description
EEH_NSTOPPED_RST_DEA	Reset deactivated and adapter is not in stopped state.
EEH_NSTOPPED_RST_ACT	Reset activated and adapter is not in stopped state.
EEH_STOPPED_LS_DIS	Adapter in stopped state with reset signal deactivated and Load/Store disabled.
EEH_STOPPED_LS_ENA	Adapter in stopped state with reset signal deactivated and Load/Store enabled.
EEH_UNAVAILABLE	Adapter is either permanently or temporarily unavailable.

Valid *support* values are as follows:

Item	Description
0	EEH not supported.
1	EEH supported.

The driver should call this service and check for `EEH_STOPPED_LS_DIS` and `EEH_STOPPED_LS_ENA` as the *state* values if it suspects an EEH event on the adapter. If the *state* is either of those values, the slot is said to be frozen.

**Single-function Driver:** A single-function adapter driver calls this service to query the state of the slot. If the service fails due to hardware or firmware reasons, an error is logged. If the service fails, *state* and *support* values are undefined, and `EEH_FAIL` is returned.

**Multifunction Driver:** For a multifunction adapter driver, this service analyzes the *state* to determine if:

- The state is frozen, or
- it is permanently unavailable (that is, the slot is unusable from hereon), or
- it is temporarily unavailable.

If the slot is in either a frozen or temporarily unavailable state, the `EEH_DD_SUSPEND` message is broadcast to all registered drivers on this slot. If the slot is permanently unavailable (that is, dead), the `EEH_DD_DEAD` message is broadcast. Upon receiving this message, the drivers are expected to suspend all further DMA, PIO, interrupt, configuration cycles, and so on until the slot is recovered. If the service fails due to hardware or firmware reasons, an error is logged, `EEH_DD_DEAD` is broadcast, and `EEH_FAIL` is returned.

### Temporarily versus permanently unavailable state

In addition to *state* and *support*, this service also returns a valid *retry\_delay* value in the `eeh_handle` structure if the *state* is `EEH_UNAVAILABLE`. If *retry\_delay* is 0, it is permanently unavailable. If *retry\_delay* is non-zero, it is temporarily unavailable. A permanently unavailable state means that the slot is unusable until a hot-plug operation or partition reboot is performed. Therefore, the drivers mark their adapters as unusable when they receive an `EEH_UNAVAILABLE` message (single-function) or when they receive an `EEH_DD_DEAD` message (multifunction). A temporarily unavailable state means that the current *state* of a slot is transient and might take a few minutes to settle down. Until that time, the device driver cannot begin recovery because it does not know what the final state will be. The temporarily unavailable state is handled differently by the single-function and multifunction drivers as follows:

**Single-function Driver:** Because a single-function driver drives its own recovery, it needs to check for *retry\_delay* if the *state* is set to `EEH_UNAVAILABLE`. If *retry\_delay* is non-zero, it represents the number of seconds that the driver should wait before calling this kernel service again. It continues to call this service repeatedly as long as the *state* is `EEH_UNAVAILABLE` and *retry\_delay* is non-zero. Eventually, the *state* will end up in one of the following:

- `EEH_NSTOPPED_RST_ACT`
- `EEH_STOPPED_LS_DIS`
- `EEH_UNAVAILABLE` w/ "retry\_delay" set to 0 (i.e. permanently unavailable)

At that point, the driver can continue with its normal course of action for a given state.

**Multifunction Driver:** A multifunction driver does not need to check for the *retry\_delay* field when the *state* is `EEH_UNAVAILABLE`, because `EEH_UNAVAILABLE` would only mean permanently unavailable. In the case of temporarily unavailable, a multifunction driver would receive the `EEH_DD_SUSPEND` or `EEH_DD_DEAD` message after some time, depending upon the final *state* of the slot. If the final state was `EEH_NSTOPPED_RST_ACT` or `EEH_STOPPED_LS_DIS`, then `EEH_DD_SUSPEND` is broadcast; if it was `EEH_UNAVAILABLE`, then `EEH_DD_DEAD` is broadcast. Thus, from the point-of-view of a multifunction driver, there is no difference between frozen and temporarily unavailable.

The macro `EEH_READ_SLOT_STATE(handle, state, support)` is provided for device drivers to call this service.

## Execution Environment

This kernel service can be called from the process or interrupt environment.

## Return Values

Item	Description
<code>EEH_SUCC</code>	Successfully read the slot state and capabilities
<code>EEH_FAIL</code>	Unable to read the slot state and capabilities

### Related reference:

“`eeh_enable_slot` Kernel Service” on page 128

“`eeh_init` Kernel Service” on page 129

“`eeh_slot_error` Kernel Service” on page 136

## eeh\_reset\_slot Kernel Service

### Purpose

This service activates, deactivates, or toggles the reset line of a PCI slot.

### Syntax

```
#include <sys/eeh.h>
```

```
long eeh_reset_slot(handle, flag)
eeh_handle_t handle;
long flag;
```

### Parameters

Item	Description
<i>handle</i>	EEH handle obtained from the <code>eeh_init</code> or <code>eeh_init_multifunc</code> kernel services
<i>flag</i>	Flag can be either <code>EEH_ACTIVE</code> or <code>EEH_DEACTIVE</code> .

### Description

**Single-function Drivers:** This service activates and deactivates the reset line between the Terminal Bridge and the adapter. The *flag* argument specifies whether to activate (`EEH_ACTIVE`) or deactivate (`EEH_DEACTIVE`) depending upon the required action. To do the reset of a slot, the reset line should be toggled by calling this service twice: once with `EEH_ACTIVE` followed by a second call with `EEH_DEACTIVE`. There should be a minimum of 100 milliseconds delay between the activation and deactivation of the signal. The minimum delay is specified by the PCI System Architecture and should be enforced by the single-function driver.

**Multifunction Drivers:** On a multifunction adapter, the `EEH_MASTER` for the slot drives error recovery. Therefore, only the `EEH_MASTER` can call this service. Unlike the single-function driver, the master calls this service only once with the `EEH_ACTIVE` flag.

For the multi-function drivers, the service first activates and then deactivates the reset signal on the slot. It enforces a 100-millisecond delay between the activation and deactivation as mandated by the PCI System Architecture. After the reset signal is deactivated, the service attempts to reconfigure the bridge on the adapter, if there is one (only applies to the bridged-adapters), after `dd_trb_timer` seconds specified in `eeh_init_multifunc()`. At the end of a successful reset and optional bridge recovery, an

EEH\_DD\_RESUME message is broadcast to the slot's multifunction drivers notifying them to resume normal operation. If this service fails, the EEH\_DD\_DEAD message is broadcast. If failure is due to hardware or firmware, an error is logged.

**EEH\_SAFE mode:** If an EEH\_SAFE adapter calls this service, the reset signal is activated but is never deactivated, thereby leaving the adapter in a "permanently unavailable" state. Such an adapter becomes available again if either the PCI hot-plug operation is performed on it or if the partition is rebooted. This service returns EEH\_FAIL for an EEH\_SAFE driver.

The macro `EEH_RESET_SLOT(handle, flag)` is provided for device drivers to call this service.

## Execution Environment

This kernel service can be called from the process or interrupt environment.

## Return Values

Item	Description
EEH_SUCC	Slot reset activate/deactivate succeeded
EEH_FAIL	Failed to activate/deactivate the reset line, nonmaster called the service, or EEH_SAFE mode is active
EEH_BUSY	Recovery is already in progress

### Related reference:

"`eeh_enable_slot` Kernel Service" on page 128

"`eeh_read_slot_state` Kernel Service" on page 133

"`eeh_slot_error` Kernel Service"

## eeh\_slot\_error Kernel Service Purpose

This service logs a temporary or permanent error and optionally marks the slot permanently unavailable.

## Syntax

```
#include <sys/eeh.h>
```

```
long eeh_slot_error(handle, flag, dd_buf, dd_buf_length)
eeh_handle_t   handle;
int           flag;
char          *dd_buf;
long         dd_buf_length;
```

## Parameters

Item	Description
<i>handle</i>	EEH handle obtained from <code>eeh_init</code> or <code>eeh_init_multifunc</code>
<i>flag</i>	EEH_RESET_TEMP or EEH_RESET_PERM
<i>dd_buf</i>	Address of the device driver's error log buffer
<i>dd_buf_length</i>	Length of device driver's error log buffer in bytes

## Description

This service performs a number of tasks:

- It collects hardware data to help in understanding the nature and source of an EEH event
- It combines the device-driver-supplied debug data log with the hardware data log and creates an entry in the error log



- It optionally marks the slot permanently unavailable so that subsequent `eeh_read_slot_state()` calls return `EEH_UNAVAILABLE` with a `retry_delay` value of 0

The behavior of this kernel service is controlled by two *flag* values:

Item	Description
<code>EEH_RESET_TEMP</code>	This flag performs only the first two of the preceding tasks.
<code>EEH_RESET_PERM</code>	This flag performs all three tasks.

Depending on the hardware state of the slot, this service might not be able to collect the hardware data. Thus, the service succeeds but logs no data. If `EEH_RESET_PERM` was supplied, it still marks the slot permanently unavailable.

The `dd_buf` and `dd_buf_length` parameters are used to combine the device driver error log with the hardware log. The `dd_buf` argument is the address of an error log buffer that contains the device driver's data. The `dd_buf_length` argument is the length of this buffer. If the length exceeds `MAX_DD_LOG_SIZE` bytes, the driver's log data is truncated. If `dd_buf` is `NULL`, the error log contains only hardware data, if any.

**Single-function driver:** The kernel service works as in the preceding description. If it fails because of hardware or firmware reasons, `EEH_FAIL` is returned and an error is logged.

**Multifunction driver:** For the multifunction drivers, this service works as in the preceding description, except that if `EEH_RESET_PERM` was supplied, the `EEH_DD_DEAD` message is broadcast.

The macro `EEH_SLOT_ERROR(handle, flag, dd_buf, dd_buf_length)` is provided for device drivers to call this service.

## Execution Environment

This kernel service can be called from the process or interrupt environment.

## Return Values

Item	Description
<code>EEH_SUCC</code>	Successfully logged error
<code>EEH_FAIL</code>	Failed to log the error and optionally mark the slot permanently unavailable

### Related reference:

“`eeh_read_slot_state` Kernel Service” on page 133

“`eeh_reset_slot` Kernel Service” on page 135

## enqueue Kernel Service

### Purpose

Sends a request queue element to a device queue.

### Syntax

```
#include <sys/types.h> #include <sys/errno.h> #include <sys/deviceq.h> int enqueue ( qe) struct req_qe
*qe;
```

### Parameter

Item	Description
<i>qe</i>	Specifies the address of the request queue element.

## Description

The **enque** kernel service is not part of the base kernel, but is provided by the device queue management kernel extension. This queue management kernel extension must be loaded into the kernel before loading any kernel extensions referencing these services.

The **enque** service places the queue element into a specified device queue. It is used for simple process-to-process communication within the kernel. The requester builds a copy of the queue element, indicated by the *qe* parameter, and passes this copy to the **enque** service. The kernel copies this queue element into a queue element in pinned global memory and then enqueues it on the target device queue.

The path identifier in the request queue element indicates the device queue into which the element is placed.

The **enque** service supports the sending of the following types of queue elements:

Queue Element	Description
SEND_CMD	Send command.
START_IO	Start I/O.
GEN_PURPOSE	General purpose.

For simple interprocess communication, general purpose queue elements are used.

The queue element priority value can range from **QE\_BEST\_PRTY** to **QE\_WORST\_PRTY**. This value is limited to the value specified when the queue was created.

The operation options in the queue element control how the queue element is processed. There are five standard operation options:

Operation Option	Description
ACK_COMPLETE	Acknowledge completion in all cases.
ACK_ERRORS	Acknowledge completion if the operation results in an error.
SYNC_REQUEST	Synchronous request.
CHAINED	Chained control blocks.
CONTROL_OPT	Kernel control operation.

**Note:** Only one of **ACK\_COMPLETE**, **ACK\_ERRORS**, or **SYNC\_REQUEST** can be specified. Also, all of these options are ignored if the path specifies that no acknowledgment (**NO\_ACK**) should be sent.

With the **SYNC\_REQUEST** synchronous request option, control does not return from the **enque** service until the request queue element is acknowledged. This performs in one step what can also be achieved by sending a queue element with the **ACK\_COMPLETE** flag on, and then calling either the **et\_wait** or **waitq** kernel services.

The kernel calls the server's **check** routine, if one is defined, before a queue element is placed on the device queue. This routine can stop the operation if it detects an error.

The kernel notifies the device queue's server, if necessary, after a queue element is placed on the device queue. This is done by posting the server process (using the **et\_post** kernel service) with an event control bit.

## Execution Environment

The **enqueue** kernel service can be called from the process environment only.

## Return Values

Item	Description
RC_GOOD	Indicates a successful operation.
RC_ID	Indicates a path identifier that is not valid.

All other error values represent errors returned by the server.

### Related reference:

“et\_post Kernel Service” on page 119

“et\_wait Kernel Service” on page 120

“waitq Kernel Service” on page 583

## errresume Kernel Service

### Purpose

Resumes error logging after an **errlast** command was issued.

### Syntax

```
void errresume()
```

### Description

When an error is logged with the **errlast** command, no more error logging will happen on the system until an **errresume** call is issued.

## Execution Environment

This can be called from either the process or an interrupt level.

### Related reference:

“errsave or errlast Kernel Service”

### Related information:

Error-Logging Facility

## errsave or errlast Kernel Service

### Purpose

Allows the kernel and kernel extensions to write to the error log.

### Syntax

```
#include <sys/types.h> #include <sys/errno.h> #include <sys/errids.h> void errsave ( buf, cnt) char *buf;  
unsigned int cnt; void errlast (buf, cnt) char *buf unsigned int cnt;
```

### Parameters

Item	Description
<i>buf</i>	Points to a buffer that contains an error record as described in the <code>/usr/include/sys/err_rec.h</code> file.
<i>cnt</i>	Specifies the number of bytes in the error record contained in the buffer pointed to by the <i>buf</i> parameter.

## Description

The **errsave** kernel service allows the kernel and kernel extensions to write error log entries to the error device driver. The error record pointed to by the *buf* parameter includes the error ID resource name and detailed data.

In addition, the **errlast** kernel service disables any future error logging, thus any error logged with **errlast** will stay on NVRAM. This service is only for use prior to a pending system crash or stop. The **errlast** service should only be used in extreme circumstances where the system can not continue, such as the occurrence of a machine check.

## Execution Environment

The **errsave** kernel service can be called from either the process or interrupt environment.

## Return Values

The **errsave** service has no return values.

### Related information:

errlog subroutine  
 Error Logging Special Files  
 RAS Kernel Services

## f

The following kernel services begin with the with the letter f.

## fetch\_and\_add Kernel Services

### Purpose

Increments a variable atomically.

### Syntax

```
#include <sys/atomic_op.h>
```

```
int fetch_and_add (addr, value)
atomic_p addr;
int value;
```

```
long fetch_and_addlp (addr, value)
atomic_l addr;
long value;
```

### Parameters

Item	Description
<i>addr</i>	Specifies the address of the variable to be incremented.
<i>value</i>	Specifies the value to be added to the variable.

## Description

The **fetch\_and\_add** kernel services atomically increment a variable.

The **fetch\_and\_add** kernel service operates on a single word (32 bit) variable while the **fetch\_and\_addlp** kernel service operates on a double word (64 bit) variable.

These operations are useful when a counter variable is shared between several kernel threads, because it ensures that the fetch, update, and store operations used to increment the counter occur atomically (are not interruptible).

### Note:

- The single word variable for the **fetch\_and\_add** kernel service must be aligned on a word (32 bit) boundary.
- The double word variable for the **fetch\_and\_addlp** kernel service must be aligned on a double word (64 bit) boundary.

## Execution Environment

The **fetch\_and\_add** kernel services can be called from either the process or interrupt environment.

## Return Values

The **fetch\_and\_add** kernel services return the original value of the variable.

### Related reference:

“**fetch\_and\_and** or **fetch\_and\_or** Kernel Services”

“**compare\_and\_swap** Kernel Services” on page 45

### Related information:

Locking Kernel Services

## **fetch\_and\_and** or **fetch\_and\_or** Kernel Services

### Purpose

Clears and sets bits in a variable atomically.

### Syntax

```
#include <sys/atomic_op.h>
```

```
uint fetch_and_and (addr, mask)
atomic_p addr;uint mask;
```

```
ulong fetch_and_andlp (addr, mask)
atomic_l addr;
ulong mask;
```

```
uint fetch_and_or (addr, mask)
atomic_p addr;
uint mask;
```

```
ulong fetch_and_orlp (addr, mask)
atomic_l addr;
ulong mask;
```

## Parameters

Item	Description
<i>addr</i>	Specifies the address of the variable whose bits are to be cleared or set.
<i>mask</i>	Specifies the bit mask which is to be applied to the variable.

## Description

The **fetch\_and\_and** and **fetch\_and\_or** kernel services respectively clear and set bits in a variable, according to a bit mask, as a single atomic operation. The **fetch\_and\_and** service clears bits in the variable which correspond to clear bits in the bit mask, and the **fetch\_and\_or** service sets bits in the variable which correspond to set bits in the bit mask.

The **fetch\_and\_add** and **fetch\_and\_or** kernel services operate on a single word (32 bit) variable while the **fetch\_and\_addlp** and **fetch\_and\_orlp** kernel services operate on a double word (64 bit) variable.

These operations are useful when a variable containing bit flags is shared between several kernel threads, because they ensure that the fetch, update, and store operations used to clear or set a bit in the variable occur atomically (are not interruptible).

### Note:

- For the **fetch\_and\_and** and **fetch\_and\_or** kernel services, the single word containing the bit flags must be aligned on a full word (32 bit) boundary.
- For the **fetch\_and\_addlp** and **fetch\_and\_orlp** kernel services, the double word containing the bit flags must be aligned on a double word (64 bit) boundary.

## Execution Environment

The **fetch\_and\_and** and **fetch\_and\_or** kernel services can be called from either the process or interrupt environment.

## Return Values

The **fetch\_and\_and** and **fetch\_and\_or** kernel services return the original value of the variable.

### Related reference:

“fetch\_and\_add Kernel Services” on page 140

“compare\_and\_swap Kernel Services” on page 45

### Related information:

Locking Kernel Services

## fidtovp Kernel Service

### Purpose

Maps a file system structure to a file ID.

Maps a file identifier to a mode.

## Syntax

```
#include <sys/types.h> #include <sys/vnode.h> int fidtovp(fsid, fid, vpp) fsid_t *fsid; struct fileid *fid;  
struct vnode **vpp;
```

## Parameters

Item	Description
<i>fsid</i>	Points to a file system ID structure. The system uses this structure to determine which virtual file system (VFS) contains the requested file.
<i>fid</i>	Points to a file ID structure. The system uses this pointer to locate the specific file within the VFS.
<i>vpp</i>	Points to a location to store the file's vnode pointer upon successful return of the <b>fidtovp</b> kernel service.

## Description

The **fidtovp** kernel service returns a pointer to a vnode for the file identified by **fsid** and **fid**, and increments the count on the vnode so the file is not removed. Subroutines that call the **fidtovp** kernel service must call VNOP\_RELE to release the vnode pointer.

This kernel service is designed for use by the server side of distributed file systems.

## Execution Environment

The **fidtovp** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates successful completion.
ESTALE	Indicates the requested file or file system was removed or recreated since last access with the given file system ID or file ID.

## find\_input\_type Kernel Service

### Purpose

Finds the given packet type in the Network Input Interface switch table and distributes the input packet according to the table entry for that type.

### Syntax

```
#include <sys/types.h> #include <sys/errno.h> #include <net/if.h> int find_input_type(type, m, ac,  
header_pointer) ushort type; struct mbuf * m; struct arpcom * ac; caddr_t header_pointer;
```

### Parameters

Item	Description
<i>type</i>	Specifies the protocol type.
<i>m</i>	Points to the <b>mbuf</b> buffer containing the packet to distribute.
<i>ac</i>	Points to the network common portion ( <b>arpcom</b> ) of the network interface on which the packet was received. This common portion is defined as follows: in <code>net/if_arp.h</code>
<i>header_pointer</i>	Points to the buffer containing the input packet header.

## Description

The **find\_input\_type** kernel service finds the given packet type in the Network Input table and distributes the input packet contained in the **mbuf** buffer pointed to by the *m* value. The *ac* parameter is passed to services that do not have a queued interface.

## Execution Environment

The **find\_input\_type** kernel service can be called from either the process or interrupt environment.

## Return Values

Item	Description
0	Indicates that the protocol type was successfully found.
ENOENT	Indicates that the service could not find the type in the Network Input table.

### Related reference:

“add\_input\_type Kernel Service” on page 8

“del\_input\_type Kernel Service” on page 65

### Related information:

Network Kernel Services

## fp\_access Kernel Service

### Purpose

Checks for access permission to an open file.

### Syntax

```
#include <sys/types.h> #include <sys/errno.h> int fp_access ( fp, perm) struct file *fp; int perm;
```

### Parameters

Item	Description
<i>fp</i>	Points to a file structure returned by the <b>fp_open</b> or <b>fp_opendev</b> kernel service.
<i>perm</i>	Indicates which read, write, and execute permissions are to be checked. The <code>/usr/include/sys/mode.h</code> file contains pertinent values (IREAD, IWRITE, IEXEC).

### Description

The **fp\_access** kernel service is used to see if either the read, write, or exec bit is set anywhere in a file's permissions mode. Set *perm* to one of the following constants from **mode.h**:

IREAD IWRITE IEXEC



## Execution Environment

The `fp_access` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates that the calling process has the requested permission.
EACCES	Indicates all other conditions.

### Related information:

access subroutine

Logical File System Kernel Services

## fp\_close Kernel Service

### Purpose

Closes a file.

### Syntax

```
#include <sys/types.h> #include <sys/errno.h> int fp_close ( fp) struct file *fp;
```

### Parameter

Item	Description
<i>fp</i>	Points to a file structure returned by the <code>fp_open</code> , <code>fp_getf</code> , or <code>fp_opendev</code> kernel service.

### Description

The `fp_close` kernel service is a common service for closing files used by both the file system and routines outside the file system.

## Execution Environment

The `fp_close` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
non-zero	The underlying file system implementation might report one of the values from the <code>/usr/include/errno.h</code> file, which is returned to the caller as a return value. However, the file is still closed.

### Related information:

close subroutine

Logical File System Kernel Services

## fp\_close Kernel Service for Data Link Control (DLC) Devices

### Purpose

Allows kernel to close the generic data link control (GDLC) device manager using a file pointer.

### Syntax

```
int fp_close( fp)
```

## Parameters

Item	Description
<i>fp</i>	Specifies the file pointer of the GDLC being closed.

## Description

The **fp\_close** kernel service disables a GDLC channel. If this is the last channel to close on a port, the GDLC device manager resets to an idle state on that port and the communications device handler is closed. The **fp\_close** kernel service may be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful completion.
ENXIO	Indicates an invalid file pointer. This value is defined in the <code>/usr/include/sys/errno.h</code> file.

### Related reference:

“fp\_close Kernel Service” on page 145

“fp\_open Kernel Service for Data Link Control (DLC) Devices” on page 157

### Related information:

Generic Data Link Control (GDLC) Environment Overview

## fp\_fstat Kernel Service

### Purpose

Gets the attributes of an open file.

### Syntax

```
#include <sys/types.h> #include <sys/errno.h> int fp_fstat (fp, statp, len, seg) struct file* fp; struct stat *statp; int len; int seg;
```

## Parameters

Item	Description
<i>fp</i>	Points to a file structure returned by the <b>fp_open</b> kernel service.
<i>statp</i>	Points to a buffer defined to be of <b>stat</b> or <b>fullstat</b> type structure. The <i>statsz</i> parameter indicates the buffer type.
<i>len</i>	Indicates the size of the <b>stat</b> or <b>fullstat</b> structure to be returned. The <code>/usr/include/sys/stat.h</code> file contains information about the <b>stat</b> structure.
<i>seg</i>	Specifies the flag indicating where the information represented by the <i>statbuf</i> parameter is located: <b>SYS_ADSPACE</b> Buffer is in kernel memory. <b>USER_ADSPACE</b> Buffer is in user memory.

## Description

The **fp\_fstat** kernel service is an internal interface to the function provided by the **fstatx** subroutine.

## Execution Environment

The `fp_fstat` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.

If an error occurs, one of the values from the `/usr/include/sys/errno.h` file is returned.

### Related information:

`fstatx` subroutine

Logical File System Kernel Services

## fp\_fsync Kernel Service Purpose

Writes changes for a specified range of a file to permanent storage.

## Syntax

```
#include <sys/fp_io.h>
```

```
int fp_fsync (fp, how, off, len)
struct file *fp;
int how;
offset_t off;
offset_t len;
```

## Description

The `fp_fsync` kernel service is an internal interface to the function provided by the `fsync_range` subroutine.

## Parameters

Item	Description
<i>fp</i>	Points to a file structure returned by the <code>fp_open</code> kernel service.
<i>how</i>	Specifies the following handling characteristics of the operation:  <b>FDATASYNC</b> The changed data in the range specified by the <code>off</code> and <code>len</code> parameters is written to the storage. If the metadata for the file is changed and this changed metadata must read the data, the metadata is also written to the storage. Otherwise, the metadata is not updated.  <b>FFILESYNC</b> The changed data in the range specified by the <code>off</code> and <code>len</code> parameters is written to the storage. If any metadata is changed, all of the changed user data is written to the storage. Metadata changes and file attributes including time stamps are also written to the storage.
<i>off</i>	Specifies the starting offset value of the data in the file to be written to the storage.
<i>len</i>	Specifies the length of the file range to be written to the storage. If you specify the value as zero, all cached data is written to the storage.

## Execution Environment

The `fp_fsync` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
ERRNO	Returns an error number from the <code>/usr/include/sys/errno.h</code> file on failure.

#### Related information:

fsync or fsync\_range Subroutine  
 Logical File System Kernel Services

### fp\_getdevno Kernel Service

#### Purpose

Gets the device number or channel number for a device.

#### Syntax

```
#include <sys/types.h> #include <sys/errno.h> #include <sys/file.h> int fp_getdevno ( fp, devp, chanp)
struct file *fp; dev_t *devp; chan_t *chanp;
```

#### Parameters

Item	Description
<i>fp</i>	Points to a file structure returned by the <code>fp_open</code> or <code>fp_opendev</code> service.
<i>devp</i>	Points to a location where the device number is to be returned.
<i>chanp</i>	Points to a location where the channel number is to be returned.

#### Description

The `fp_getdevno` service finds the device number and channel number for an open device that is associated with the file pointer specified by the *fp* parameter. If the value of either *devp* or *chanp* parameter is null, this service does not attempt to return any value for the argument.

#### Execution Environment

The `fp_getdevno` kernel service can be called from the process environment only.

#### Return Values

Item	Description
0	Indicates a successful operation.
EINVAL	Indicates that the pointer specified by the <i>fp</i> parameter does not point to a file structure for an open device.

#### Related information:

Logical File System Kernel Services

### fp\_getea Kernel Service

#### Purpose

Reads the value of an extended attribute value.

#### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
int fp_getea (fp, name, value, size, countp, segflag)
struct file * fp;
const char * name;
```

```
void * value;
size_t size;
ssize_t * countp;
int segflag;
```

## Parameters

Item	Description
<i>fp</i>	Specifies the file structure returned by the <b>fp_open</b> kernel service.
<i>name</i>	Specifies the name of the extended attribute. An extended attribute name is a NULL-terminated string.
<i>value</i>	Specifies the pointer to a buffer in which the attribute is stored. The value of an extended attribute is an opaque byte stream of specified length.
<i>size</i>	Specifies the size of the value buffer. If size is 0, <b>fp_getea</b> returns the current size of the named extended attribute, which can be used to estimate whether the size of a buffer is sufficiently large to hold the value associated with the extended attribute.
<i>countp</i>	Specifies the actual size of the content in the value buffer.
<i>segflag</i>	Specifies the flag indicating where the pointer specified by the path parameter is located: <b>SYS_ADSPACE</b> The pointers specified by the <i>name</i> and <i>value</i> parameters are stored in kernel memory. <b>USER_ADSPACE</b> The pointers specified by the <i>name</i> and <i>value</i> parameters are stored in application memory.

## Description

The **fp\_getea** kernel service provides a common service used by:

- The file system for the implementation of the **fgetea** subroutine.
- Kernel routines outside the file system that set extended attribute values.

## Execution Environment

The **fp\_getea** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
ERRNO	Indicates a failed operation. Returns an error number from the <code>/usr/include/sys/errno.h</code> file on failure.

### Related information:

fgetea subroutine

fp\_open subroutine

Logical File System Kernel Services

## fp\_getf Kernel Service

### Purpose

Retrieves a pointer to a file structure.

### Syntax

```
#include <sys/types.h> #include <sys/errno.h> int fp_getf ( fd, fpp) int fd; struct file **fpp;
```

## Parameters

Item	Description
<i>fd</i>	Specifies a file descriptor.
<i>fp</i>	Points to the location where the file pointer is to be returned.

## Description

A process calls the **fp\_getf** kernel service when it has a file descriptor for an open file, but needs a file pointer to use other Logical File System services.

The **fp\_getf** kernel service uses the file descriptor as an index into the process's open file table. From this table it extracts a pointer to the associated file structure.

As a side effect of the call to the **fp\_getf** kernel service, the reference count on the file descriptor is incremented. This count must be decremented when the caller has completed its use of the returned file pointer. The file descriptor reference count is decremented by a call to the **ufdrele** kernel service.

## Execution Environment

The **fp\_getf** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
EBADF	Indicates that either the file descriptor is invalid or not currently used in the process.

### Related reference:

“ufdhold and ufdrele Kernel Service” on page 516

### Related information:

Logical File System Kernel Services

## fp\_get\_path Kernel Service

### Purpose

Returns the full path name of the file referenced by the **fp** parameter.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>

int
fp_get_path(struct file *fp,
            int flags,
            char *path,
            size_t size)
```

## Parameters

**fp** Points to a file structure that is returned by the **fp\_open** or **fp\_opendev** kernel service.

### flags

No flags are defined; this parameter must be 0.

### path

Points to a buffer where the file name is returned.

**size**

Specifies the size of the path buffer.

**Description**

The `fp_get_path` kernel service provides a method to find a path name from a file structure pointer.

**Execution environment**

The `fp_get_path` kernel service can be called only from the process environment.

**Return values**

0 Indicates a successful operation.

**EINVAL**

Invalid `fp` or `path` argument, or the `fp` parameter does not refer to a `DTYPE_VNODE` file structure.

**fp\_hold Kernel Service****Purpose**

Increments the open count for a specified file pointer.

**Syntax**

```
#include <sys/types.h>
```

```
#include <sys/errno.h>
```

```
void fp_hold ( fp)
```

```
struct file *fp;
```

**Parameter**

Item	Description
------	-------------

<i>fp</i>	Points to a file structure previously obtained by calling the <code>fp_open</code> , <code>fp_getf</code> , or <code>fp_opendev</code> kernel service.
-----------	--

**Description**

The `fp_hold` kernel service increments the use count in the file structure specified by the `fp` parameter. This results in the associated file remaining opened even when the original open is closed.

If this function is used, and access to the file associated with the pointer specified by the `fp` parameter is no longer required, the `fp_close` kernel service should be called to decrement the use count and close the file as required.

**Execution Environment**

The `fp_hold` kernel service can be called from the process environment only.

**Related information:**

Logical File System Kernel Services

**fp\_ioctl Kernel Service****Purpose**

Issues a control command to an open device or file.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
int fp_ioctl (fp, cmd, arg, ext)
struct file * fp;
unsigned long cmd;
caddr_t arg;
int ext;
```

## Parameters

Item	Description
<i>fp</i>	Points to a file structure returned by the <b>fp_open</b> or <b>fp_opendev</b> kernel service.
<i>cmd</i>	Specifies the specific control command requested.
<i>arg</i>	Indicates the data required for the command.
<i>ext</i>	Specifies an extension argument required by some device drivers. Its content, form, and use are determined by the individual driver.

## Description

The **fp\_ioctl** kernel service is an internal interface to the function provided by the **ioctl** subroutine.

## Execution Environment

The **fp\_ioctl** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.

If an error occurs, one of the values from the **/usr/include/sys/errno.h** file is returned. The **ioctl** subroutine contains valid **errno** values.

### Related information:

[ioctl subroutine](#)

[Logical File System Kernel Services](#)

## **fp\_ioctl** Kernel Service for Data Link Control (DLC) Devices

### Purpose

Transfers special commands from the kernel to generic data link control (GDLC) using a file pointer.

## Syntax

```
#include <sys/gdlextc.h>
#include <fcntl.h>
```

```
int fp_ioctl (fp, cmd, arg, ext)
```

## Parameters



Item	Description
<i>fp</i>	Specifies the file pointer of the target GDLC.
<i>cmd</i>	Specifies the operation to be performed by GDLC. For a listing of all possible operators, see "ioctl Operations (op) for DLC" <i>Technical Reference: Communications, Volume 1</i> .
<i>arg</i>	Specifies the address of the parameter block. The argument for this parameter must be in the kernel space. For a listing of possible values, see "Parameter Blocks by ioctl Operation for DLC" <i>Technical Reference: Communications, Volume 1</i> .
<i>ext</i>	Specifies the extension parameter. This parameter is ignored by GDLC.

## Description

Various GDLC functions can be initiated using the **fp\_ioctl** kernel service, such as changing configuration parameters, contacting the remote, and testing a link. Most of these operations can be completed before returning to the user synchronously. Some operations take longer, so asynchronous results are returned much later using the **exception** function handler. GDLC calls the kernel user's exception handler to complete these results. Each GDLC supports the **fp\_ioctl** kernel service by way of its **dlcioctl** entry point. The **fp\_ioctl** kernel service may be called from the process environment only.

**Note:** The **DLC\_GET\_EXCEP** ioctl operation is not used since all exception conditions are passed to the kernel user through the exception handler.

## Return Values

Item	Description
0	Indicates a successful completion.
ENXIO	Indicates an invalid file pointer.
EINVAL	Indicates an invalid value.
ENOMEM	Indicates insufficient resources to satisfy the <b>ioctl</b> subroutine.

These return values are defined in the `/usr/include/sys/errno.h` file.

### Related reference:

"fp\_ioctl Kernel Service" on page 151

### Related information:

ioctl subroutine

Generic Data Link Control (GDLC) Environment Overview

## fp\_ioctlx Kernel Service

### Purpose

Issues a control command to an open device.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <fcntl.h>
```

```
int fp_ioctlx (fp, cmd, arg, ext, flags, retval)
struct file *fp;
unsigned long cmd;
caddr_t arg;
ext_t ext;
unsigned long flags;
long *retval;
```

## Description

The `fp_ioctlx` kernel service is an internal interface to the function provided by the `ioctl` subroutine.

The `fp_ioctlx` kernel service issues a control command to an open device. Some drivers need the return value that is returned by the kernel service if there is no error. This value is not available through the `fp_ioctl` kernel service. The `fp_ioctlx` kernel service allows this data to be passed.

## Parameters

Item	Description
<i>fp</i>	Points to a file structure returned by the <code>fp_open</code> or <code>fp_opendev</code> kernel service.
<i>cmd</i>	Specifies the specific control command requested.
<i>arg</i>	Indicates the data required for the command.
<i>ext</i>	Specifies an extension argument required by some device drivers. Its content, form, and use are determined by the individual driver.
<i>flags</i>	Indicates the address space of <i>arg</i> parameter. If the <i>arg</i> value is in kernel address space, <i>flags</i> should be specified as <code>FKERNEL</code> . Otherwise, it should be zero (drivers pass data that is in user space).
<i>retval</i>	Points to the location where the return value will be stored on successful return from the call.

## Execution Environment

The `fp_ioctlx` kernel service can be called only from the process environment.

## Return Values

Upon successful completion, the `fp_ioctlx` kernel service returns 0. If unsuccessful, one of the values from the `/usr/include/sys/errno.h` file is returned. The `ioctl` subroutine contains valid `errno` values. This value will be stored in the *retval* parameter.

### Related reference:

“`fp_ioctl` Kernel Service” on page 151

### Related information:

`ioctl`, `ioctlx`, `ioctl32`, or `ioctl32x` Subroutine

## `fp_listea` Kernel Service

### Purpose

Lists the extended attributes associated with a file.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
int fp_listea (fp, list, size, countp, segflag)
struct file * fp;
const char * list;
size_t size;
ssize_t * countp;
int segflag;
```

### Parameters

<b>Item</b>	<b>Description</b>
<i>fp</i>	Specifies the file structure returned by the <b>fp_open</b> kernel service.
<i>list</i>	Specifies a pointer to a buffer in which the list of attributes will be stored.
<i>size</i>	Specifies the size of the list buffer.
<i>countp</i>	Specifies the actual size of the content in the list buffer.
<i>segflag</i>	Specifies the flags indicating where the pointer specified by the path parameter is located:
	<b>SYS_ADSPACE</b> The pointer specified by the list parameter is stored in kernel memory.
	<b>USER_ADSPACE</b> The pointer specified by the list parameter is stored in application memory.

## Description

The **fp\_listea** kernel service provides a common service used by:

- File system for the implementation of the **flistea** subroutine.
- Kernel routines outside the file system that set extend attribute values.

## Execution Environment

The **fp\_listea** kernel service can be called from the process environment only.

## Return Values

<b>Item</b>	<b>Description</b>
0	Indicates a successful operation.
ERRNO	Indicates a failed operation. Returns an error number from the <code>/usr/include/sys/errno.h</code> file on failure.

### Related information:

`flistea` subroutine

`fp_open` subroutine

Logical File System Kernel Services

## **fp\_lseek, fp\_llseek** Kernel Service Purpose

Changes the current offset in an open file.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
int fp_lseek ( fp, offset, whence)
struct file *fp;
off_t offset;
int whence;
```

```
int fp_llseek
( fp, offset, whence)
struct file *fp
offset_t offset;
int whence;
```

## Parameters

Item	Description
<i>fp</i>	Points to a file structure returned by the <b>fp_open</b> kernel service.
<i>offset</i>	Specifies the number of bytes (positive or negative) to move the file pointer.
<i>whence</i>	Indicates how to use the offset value:  <b>SEEK_SET</b> Sets file pointer equal to the number of bytes specified by the <i>offset</i> parameter.  <b>SEEK_CUR</b> Adds the number of bytes specified by the <i>offset</i> parameter to current file pointer.  <b>SEEK_END</b> Adds the number of bytes specified by the <i>offset</i> parameter to current end of file.

## Description

The **fp\_lseek** and **fp\_llseek** kernel services are internal interfaces to the function provided by the **lseek** and **llseek** subroutines.

## Execution Environment

The **fp\_lseek** and **fp\_llseek** kernel services can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
ERRNO	Returns an error number from the <code>/usr/include/sys/errno.h</code> file on failure.

### Related information:

`lseek` subroutine

Logical File System Kernel Services

## fp\_open Kernel Service

### Purpose

Opens special and regular files or directories.

### Syntax

```
#include <sys/types.h>  
#include <sys/errno.h>
```

```
int fp_open (path, oflags, mode, ext, segflag, fpp)  
char * path;  
long oflags;  
int mode;  
ext_t ext;  
int segflag;  
struct file ** fpp;
```

## Parameters

<b>Item</b>	<b>Description</b>
<i>path</i>	Points to the file name of the file to be opened.
<i>oflags</i>	Specifies open mode flags as described in the <b>open</b> subroutine.
<i>mode</i>	Specifies the mode (permissions) value to be given to the file if the file is to be created.
<i>ext</i>	Specifies an extension argument required by some device drivers. Individual drivers determine its content, form, and use.
<i>segflag</i>	Specifies the flag indicating where the pointer specified by the <i>path</i> parameter is located:  <b>SYS_ADSPACE</b> The pointer specified by the <i>path</i> parameter is stored in kernel memory.  <b>USER_ADSPACE</b> The pointer specified by the <i>path</i> parameter is stored in application memory.
<i>fpp</i>	Points to the location where the file structure pointer is to be returned by the <b>fp_open</b> service.

## Description

The **fp\_open** kernel service provides a common service used by:

- The file system for the implementation of the **open** subroutine
- Kernel routines outside the file system that must open files

## Execution Environment

The **fp\_open** kernel service can be called from the process environment only.

## Return Values

<b>Item</b>	<b>Description</b>
0	Indicates a successful operation.

Also, the *fpp* parameter points to an open file structure that is valid for use with the other Logical File System services. If an error occurs, one of the values from the **/usr/include/sys/errno.h** file is returned. The discussion of the **open** subroutine contains possible **errno** values.

### Related information:

open subroutine

Logical File System Kernel Services

## **fp\_open** Kernel Service for Data Link Control (DLC) Devices Purpose

Allows kernel to open the generic data link control (GDLC) device manager by its device name.

### Syntax

```
#include <sys/gd1extcb.h>
```

```
#include <fcntl.h>
```

```
fp_open (path, oflags, cmode, ext, segflag, fpp)
```

### Parameters

<b>Item</b>	<b>Description</b>
<i>path</i>	Consists of a character string containing the <i>/dev</i> special file name of the GDLC device manager, with the name of the communications device handler appended. The format is shown in the following example: <i>/dev/dlccether/ent0</i>
<i>oflags</i>	Specifies a value to set the file status flag. The GDLC device manager ignores all but the following values: <b>O_RDWR</b> Open for reading and writing. This must be set for GDLC or the open will not be successful. <b>O_NDELAY, O_NONBLOCK</b> Subsequent writes return immediately if no resources are available. The calling process is not put to sleep.
<i>cmode</i>	Specifies the <b>O_CREAT</b> mode parameter. This is ignored by GDLC.
<i>ext</i>	Specifies the extended kernel service parameter. This is a pointer to the <b>dlc_open_ext</b> extended I/O structure for open subroutines. The argument for this parameter must be in the kernel space. "open Subroutine Extended Parameters for DLC" <i>Technical Reference: Communications, Volume 1</i> provides more information on the extension parameter.
<i>segflag</i>	Specifies the segment flag indicating where the <i>path</i> parameter is located: <b>FP_SYS</b> The <i>path</i> parameter is stored in kernel memory.
<i>fpp</i>	<b>FP_USR</b> The <i>path</i> parameter is stored in application memory. Specifies the returned file pointer. This parameter is passed by reference and updated by the file I/O subsystem to be the file pointer for this <b>open</b> subroutine.

## Description

The **fp\_open** kernel service allows the kernel user to open a GDLC device manager by specifying the special file names of both the DLC and the communications device handler. Since the GDLC device manager is multiplexed, more than one process can open it (or the same process multiple times) and still have unique channel identifications.

Each open carries the communications device handler's special file name so that the DLC knows which port to transfer data on.

The kernel user must also provide functional entry addresses in order to obtain receive data and exception conditions. Each GDLC supports the **fp\_open** kernel service via its **dlcopen** entry point. The **fp\_open** kernel service may be called from the process environment only. "Using GDLC Special Kernel Services" in *AIX Version 6.1 Communications Programming Concepts* provides additional information.

## Return Values

Upon successful completion, this service returns a value of 0 and a valid file pointer in the *fpp* parameter.

Item	Description
ECHILD	Indicates that the service cannot create a kernel process.
EINVAL	Indicates an invalid value.
ENODEV	Indicates that no such device handler is present.
ENOMEM	Indicates insufficient resources to satisfy the open.
EFAULT	Indicates that the kernel service, such as the <b>copyin</b> or <b>initp</b> service, has failed.

These return values are defined in the `/usr/include/sys/errno.h` file.

#### Related reference:

“fp\_open Kernel Service” on page 156

“fp\_close Kernel Service for Data Link Control (DLC) Devices” on page 145

#### Related information:

open Subroutine Extended Parameters for DLC

Generic Data Link Control (GDLC) Environment Overview

## fp\_opendev Kernel Service

### Purpose

Opens a device special file.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
int fp_opendev (devno, flags, cname, ext, fpp)
dev_t devno;
int flags;
caddr_t cname;
ext_t ext;
struct file** fpp;
```

### Parameters

Item	Description
<i>devno</i>	Specifies the major and minor device number of device driver to open.
<i>flag</i>	Specifies one of the following values:  <b>DREAD</b> The device is being opened for reading only.  <b>DWRITE</b> The device is being opened for writing.  <b>DNDELAY</b> The device is being opened in nonblocking mode.
<i>cname</i>	Points to a channel specifying a character string or a null value.
<i>ext</i>	Specifies an extension argument required by some device drivers. Its content, form, and use are determined by the individual driver.
<i>fpp</i>	Specifies the returned file pointer. This parameter is passed by reference and is updated by the <b>fp_opendev</b> service to be the file pointer for this open instance. This file pointer is used as input to other Logical File System services to specify the open instance.

## Description

The kernel or kernel extension calls the **fp\_opendev** kernel service to open a device by specifying its device major and minor number. The **fp\_opendev** kernel service provides the correct semantics for opening the character or multiplexed class of device drivers.

If the specified device driver is non-multiplexed:

- An in-core i-node is found or created for this device.
- The i-node reference count is incremented by 1.
- The device driver's **ddopen** entry point is called with the *devno*, *devflag*, and *ext* parameters. The unused *chan* parameter on the call to the **ddopen** routine is set to 0.

If the device driver is a multiplexed character device driver (that is, its **ddmpx** entry point is defined), an in-core i-node is created for this channel. The device driver's **ddmpx** routine is also called with the *cname* pointer to the channel identification string if non-null. If the *cname* pointer is null, the **ddmpx** device driver routine is called with the pointer to a null character string.

If the device driver can allocate the channel, the **ddmpx** routine returns a channel ID, represented by the *chan* parameter. If the device driver cannot allocate a channel, the **fp\_opendev** kernel service returns an **ENXIO** error code. If successful, the i-node reference count is incremented by 1. The device driver **ddopen** routine is also called with the *devno*, *flag*, *chan* (provided by **ddmpx** routine), and *ext* parameters.

If the return value from the specified device driver **ddopen** routine is nonzero, it is returned as the return code for the **fp\_opendev** kernel service. If the return code from the device driver **ddopen** routine is 0, the **fp\_opendev** service returns the file pointer corresponding to this open of the device.

The **fp\_opendev** kernel service can only be called in the process environment or device driver top half. Interrupt handlers cannot call it. It is assumed that all arguments to the **fp\_opendev** kernel service are in kernel space.

The file pointer (*fpp*) returned by the **fp\_opendev** kernel service is only valid for use with a subset of the Logical File System services. These nine services can be called:

- **fp\_close**
- **fp\_ioctl**
- **fp\_poll**
- **fp\_select**
- **fp\_read**
- **fp\_readv**
- **fp\_rwuio**
- **fp\_write**
- **fp\_writv**

Other services return an **EINVAL** return value if called.

## Execution Environment

The **fp\_opendev** kernel service can be called from the process environment only.

## Return Values



Item	Description
0	Indicates a successful operation.

The *\*fpp* field also points to an open file structure that is valid for use with the other Logical File System services. If an error occurs, one of the following values from the `/usr/include/sys/errno.h` file is returned:

Item	Description
EINVAL	Indicates that the major portion of the <i>devno</i> parameter exceeds the maximum number allowed, or the <i>flags</i> parameter is not valid.
ENODEV	Indicates that the device does not exist.
EINTR	Indicates that the signal was caught while processing the <code>fp_opendev</code> request.
ENFILE	Indicates that the system file table is full.
ENXIO	Indicates that the device is multiplexed and unable to allocate the channel.

The `fp_opendev` service also returns any nonzero return code returned from a device driver `ddopen` routine.

**Related reference:**

“`ddopen` Device Driver Entry Point” on page 629

“`fp_close` Kernel Service” on page 145

**Related information:**

Logical File System Kernel Services

## fp\_poll Kernel Service

### Purpose

Checks the I/O status of multiple file pointers, file descriptors, and message queues.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/poll.h>
```

```
int fp_poll (listptr, nfdsmgs, timeout, flags)
void * listptr;
unsigned long nfdsmgs;
long timeout;
uint flags;
```

### Parameters

Item	Description
<i>listptr</i>	Points to an array of <code>pollfd</code> or <code>pollmsg</code> structures, or to a single <code>pollist</code> structure. Each structure specifies a file pointer, file descriptor, or message queue ID. The events of interest for this file or message queue are also specified.
<i>nfdsmgs</i>	Specifies the number of files and message queues to check. The low-order 16 bits give the number of elements present in the array of <code>pollfd</code> structures. The high-order 16 bits give the number of elements present in the array of <code>pollmsg</code> structures. If either half of the <i>nfdsmgs</i> parameter is equal to 0, then the corresponding array is presumed <code>abse1e</code> .
<i>timeout</i>	Specifies how long the service waits for a specified event to occur. If the value of this parameter is -1, the <code>fp_poll</code> kernel service does not return until at least one of the specified events has occurred. If the time-out value is 0, the <code>fp_poll</code> kernel service does not wait for an event to occur. Instead, the service returns immediately even if none of the specified events have occurred. For any other value of the <i>timeout</i> parameter, the <code>fp_poll</code> kernel service specifies the maximum length of time (in milliseconds) to wait for at least one of the specified events to occur.

Item	Description
<i>flags</i>	Specifies the type of data in the <i>listptr</i> parameter:
<b>POLL_FDMMSG</b>	Input is a file descriptor and/or message queue.
0	Input is a file pointer.

## Description

**Note:** The **fp\_poll** service applies only to character devices, pipes, message queues, and sockets. Not all character device drivers support the **fp\_poll** service.

The **fp\_poll** kernel service checks the specified file pointers/descriptors and message queues to see if they are ready for reading or writing, or if they have an exceptional condition pending.

The **pollfd**, **pollmsg**, and **pollist** structures are defined in the **/usr/include/sys/poll.h** file. These are the same structures described for the **poll** subroutine. One difference is that the **fd** field in the **pollfd** structure contains a file pointer when the *flags* parameter on the **fp\_poll** kernel service equals 0 (zero). If the *flags* parameter is set to a **POLL\_FDMMSG** value, the field is taken as a file descriptor in all processed **pollfd** structures. If either the **fd** or **msgid** fields in their respective structures has a negative value, the processing for that structure is skipped.

When performing a poll operation on both files and message queues, the *listptr* parameter points to a **pollist** structure, which can specify both files and message queues. To construct a **pollist** structure, use the **POLLIST** macro as described in the **poll** subroutine.

If the number of **pollfd** elements in the *nfdsmgs* parameter is 0, then the *listptr* parameter must point to an array of **pollmsg** structures.

If the number of **pollmsg** elements in the *nfdsmgs* parameter is 0, then the *listptr* parameter must point to an array of **pollfd** structures.

If the number of **pollmsg** and **pollfd** elements are both nonzero in the *nfdsmgs* parameter, the *listptr* parameter must point to a **pollist** structure as previously defined.

## Execution Environment

The **fp\_poll** kernel service can be called from the process environment only.

## Return Values

Upon successful completion, the **fp\_poll** kernel service returns a value that indicates the total number of files and message queues that satisfy the selection criteria. The return value is similar to the *nfdsmgs* parameter in the following ways:

- The low-order 16 bits give the number of files.
- The high-order 16 bits give the number of message queue identifiers that have nonzero *revents* values.

Use the **NFDS** and **NMSGs** macros to separate these two values from the return value. A return code of 0 (zero) indicates that:

- The call has timed out.
- None of the specified files or message queues indicates the presence of an event.

In other words, all *revents* fields are 0 (zero).

When the return code from the **fp\_poll** kernel service is negative, it is set to the following value:

Item	Description
EINTR	Indicates that a signal was caught during the <code>fp_poll</code> kernel service.

#### Related reference:

“selreg Kernel Service” on page 465

#### Related information:

poll subroutine

Logical File System Kernel Services

## fp\_read Kernel Service

### Purpose

Performs a read on an open file with arguments passed.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
int fp_read (fp, buf, nbytes, ext, segflag, countp)
struct file * fp;
char * buf;
ssize_t nbytes;
ext_t ext;
int segflag;
ssize_t * countp;
```

### Parameters

Item	Description
<i>fp</i>	Points to a file structure returned by the <code>fp_open</code> or <code>fp_opendev</code> kernel service.
<i>buf</i>	Points to the buffer where data read from the file is to be stored.
<i>nbytes</i>	Specifies the number of bytes to be read from the file into the buffer.
<i>ext</i>	Specifies an extension argument required by some device drivers. Its content, form, and use are determined by the individual driver.
<i>segflag</i>	Indicates in which part of memory the buffer specified by the <i>buf</i> parameter is located: <ul style="list-style-type: none"> <li><b>SYS_ADSPACE</b> The buffer specified by the <i>buf</i> parameter is in kernel memory.</li> <li><b>USER_ADSPACE</b> The buffer specified by the <i>buf</i> parameter is in application memory.</li> </ul>
<i>countp</i>	Points to the location where the count of bytes actually read from the file is to be returned.

### Description

The `fp_read` kernel service is an internal interface to the function provided by the `read` subroutine.

### Execution Environment

The `fp_read` kernel service can be called from the process environment only.

### Return Values

Item	Description
0	Indicates successful completion.

If an error occurs, one of the values from the `/usr/include/sys/errno.h` file is returned.

#### Related information:

read subroutine

Logical File System Kernel Services

## fp\_readv Kernel Service

### Purpose

Performs a read operation on an open file with arguments passed in `iovec` elements.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
int fp_readv
(fp, iov, iovcnt, ext,
seg, countp)
struct file * fp;
struct iovec * iov;
ssize_t iovcnt;
ext_t ext;
int seg;
ssize_t countp;
```

### Parameters

Item	Description
<i>fp</i>	Points to a file structure returned by the <code>fp_open</code> kernel service.
<i>iov</i>	Points to an array of <code>iovec</code> elements. Each <code>iovec</code> element describes a buffer where data to be read from the file is to be stored.
<i>iovcnt</i>	Specifies the number of <code>iovec</code> elements in the array pointed to by the <i>iov</i> parameter.
<i>ext</i>	Specifies an extension argument required by some device drivers. Its content, form, and use are determined by the individual driver.
<i>seg</i>	Indicates in which part of memory the array specified by the <i>iov</i> parameter is located: <ul style="list-style-type: none"> <li><b>SYS_ADSPACE</b> The array specified by the <i>iov</i> parameter is in kernel memory.</li> <li><b>USER_ADSPACE</b> The array specified by the <i>iov</i> parameter is in application memory.</li> </ul>
<i>countp</i>	Points to the location where the count of bytes actually read from the file is to be returned.

### Description

The `fp_readv` kernel service is an internal interface to the function provided by the `readv` subroutine.

### Execution Environment

The `fp_readv` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.

If an error occurs, one of the values from the `/usr/include/sys/errno.h` file is returned.

### Related information:

`readv` subroutine

Logical File System Kernel Services

## fp\_removeea Kernel Service Purpose

Removes an extended attribute.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
int fp_removeea (fp, name, segflag)
struct file * fp;
const char * name;
int segflag;
```

### Parameters

Item	Description
<i>fp</i>	Specifies a file structure returned by the <code>fp_open</code> kernel service.
<i>name</i>	Specifies the name of the extended attribute. An extended attribute name is a NULL-terminated string.
<i>segflag</i>	Specifies the flag indicating where the pointer specified by the <i>name</i> parameter is located: <b>SYS_ADSPACE</b> The pointer specified by the <i>name</i> parameter is stored in kernel memory. <b>USER_ADSPACE</b> The pointer specified by the <i>name</i> parameter is stored in application memory.

### Description

The `fp_removeea` kernel service provides a common service used by:

- The file system for the implementation of the `fremoveea` subroutine
- Kernel routines outside the file system that set extended attribute values

### Execution Environment

The `fp_removeea` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
ERRNO	Indicates a failed operation. Returns an error number from the <code>/usr/include/sys/errno.h</code> file on failure.

#### Related information:

removeea subroutine

fp\_open subroutine

Logical File System Kernel Services

## fp\_rwuio Kernel Service

### Purpose

Performs read and write on an open file with arguments passed in a **uio** structure.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
int fp_rwuio
( fp, rw, uiop, ext)
struct file *fp;
enum uio_rw rw;
struct uio *uiop;
ext_t ext;
```

### Parameters

Item	Description
<i>fp</i>	Points to a file structure returned by the <b>fp_open</b> or <b>fp_opendev</b> kernel service.
<i>rw</i>	Indicates whether this is a read operation or a write operation. It has a value of <b>UIO_READ</b> or <b>UIO_WRITE</b> .
<i>uiop</i>	Points to a <b>uio</b> structure, which contains information such as where to move data and how much to move.
<i>ext</i>	Specifies an extension argument required by some device drivers. Its content, form, and use are determined by the individual driver.

### Description

The **fp\_rwuio** kernel service is not the preferred interface for read and write operations. The **fp\_rwuio** kernel service should only be used if the calling routine has been passed a **uio** structure. If the calling routine has not been passed a **uio** structure, it should not attempt to construct one and call the **fp\_rwuio** kernel service with it. Rather, it should pass the requisite **uio** components to the **fp\_read** or **fp\_write** kernel services.

### Execution Environment

The **fp\_rwuio** kernel service can be called from the process environment only.

### Return Values

Item	Description
0	Indicates a successful operation.

If an error occurs, one of the values from the `/usr/include/sys/errno.h` file is returned.

**Related reference:**

“`uio` Structure” on page 639

**Related information:**

Logical File System Kernel Services

## **fp\_select Kernel Service**

### **Purpose**

Provides for cascaded, or redirected, support of the `select` or `poll` request.

### **Syntax**

```
#include <sys/types.h> #include <sys/errno.h> int fp_select ( fp, events, rtneventp, notify) struct file *fp;
ushort events; ushort *rtneventp; void (*notify)();
```

### **Parameters**

Item	Description
<i>fp</i>	Points to the open instance of the device driver, socket, or pipe for which the low-level <code>select</code> operation is intended.
<i>events</i>	Identifies the events that are to be checked. There are three standard event flags defined for the <code>poll</code> and <code>select</code> functions and one informational flag. The <code>/usr/include/sys/poll.h</code> file details the event bit definition. The four basic indicators are:  <b>POLLIN</b> Input is present for the specified object.  <b>POLLOUT</b> The specified file object is capable of accepting output.  <b>POLLPRI</b> An exception condition has occurred on the specified object.  <b>POLLSYNC</b> This is a synchronous request only. If none of the requested events are true, the selected routine should not remember this request as pending. That is, the routine does not need to call the <code>selnotify</code> service because of this request.
<i>rtneventp</i>	Indicates the returned events pointer. This parameter, passed by reference, is used to indicate which selected events are true at the current time. The returned event bits include the requested events plus an additional error event indicator:  <b>POLLERR</b> An error condition was indicated by the object's <code>select</code> routine. If this flag is set, the nonzero return code from the specified object's <code>select</code> routine is returned as the return code from the <code>fp_select</code> kernel service.
<i>notify</i>	Points to a routine to be called when the specified object invokes the <code>selnotify</code> kernel service for an outstanding asynchronous <code>select</code> or <code>poll</code> event request. If no routine is to be called, this parameter must be <code>NULL</code> .

### **Description**

The `fp_select` kernel service is a low-level service used by kernel extensions to perform a `select` operation for an open device, socket, or named pipe. The `fp_select` kernel service can be used for both synchronous and asynchronous `select` requests. Synchronous requests report on the current state of a device, and asynchronous requests allow the caller to be notified of future events on a device.

## Invocation from a Device Driver's `ddselect` Routine

A device driver's `ddselect` routine can call the `fp_select` kernel service to pass select/poll requests to other device drivers. The `ddselect` routine for one device invokes the `fp_select` kernel service, which calls the `ddselect` routine for a second device, and so on. This is required when event information for the original device depends upon events occurring on other devices. A cascaded chain of select requests can be initiated that involves more than two devices, or a single device can issue `fp_select` calls to several other devices.

Each `ddselect` routine should preserve, in its call to the `fp_select` kernel service, the same `POLLSYNC` indicator that it received when previously called by the `fp_select` kernel service.

## Invocation from Outside a Device Driver's `ddselect` Routine

If the `fp_select` kernel service is invoked outside of the device driver's `ddselect` routine, the `fp_select` kernel service sets the `POLLSYNC` flag, always making the request synchronous. In this case, no notification of future events for the specified device occurs, nor is a `notify` routine called, if specified. The `fp_select` kernel service can be used in this manner (unrelated to a poll or select request in progress) to check an object's current status.

## Asynchronous Processing and the Use of the `notify` Routine

For asynchronous requests, the `fp_select` kernel service allows its callers to register a `notify` routine to be called by the kernel when specified events become true. When the relevant device driver detects that one or more pending events have become true, it invokes the `selnotify` kernel service. The `selnotify` kernel service then calls the `notify` routine, if one has been registered. Thus, the `notify` routine is called at interrupt time and must be programmed to run in an interrupt environment.

Use of a `notify` routine affects both the calling sequence at interrupt time and how the requested information is actually reported. Generalized asynchronous processing entails the following sequence of events:

1. A select request is initiated on a device and passed on (by multiple `fp_select` kernel service invocations) to further devices. Eventually, a device driver's `ddselect` routine that is not dependent on other devices for information is reached. This `ddselect` routine finds that none of the requested events are true, but remembers the asynchronous request, and returns to the caller. In this way, the entire chain of calls is backed out, until the origin of the select request is reached. The kernel then puts the originating process to sleep.
2. Later, one or more events become true for the device remembering the asynchronous request. The device driver routine (possibly an interrupt handler) calls the `selnotify` kernel service.
3. If the events are still being waited on, the `selnotify` kernel service responds in one of two ways. If no `notify` routine was registered when the select request was made for the device, then all processes waiting for events on this device are awakened. If a `notify` routine exists for the device, then this routine is called. The `notify` routine determines whether the original requested event should be reported as true, and if so, calls the `selnotify` kernel service on its own.

The following example details a cascaded scenario involving several devices. Suppose that a request has been made for Device A, and Device A depends on Device B, which depends on Device C. When specified events become true at Device C, the `selnotify` kernel service called from Device C's device driver performs differently depending on whether a `notify` routine was registered at the time of the request.

## Cascaded Processing without the Use of `notify` Routines

If no `notify` routine was registered from Device B, then the `selnotify` kernel service determines that the specified events are to be considered true for the device driver at the head of the cascading chain. (The



head of the chain, in this case Device A, is the first device driver to issue the **fp\_select** kernel service from its select routine.) The **selnotify** kernel service awakens all processes waiting for events that have occurred on Device A.

It is important to note that when no **notify** routine is used, any device driver in the calling chain that reports an event with the **selnotify** kernel service causes that event to appear true for the first device in the chain. As a result, any processes waiting for events that have occurred on that first device are awakened.

### Cascaded Processing with notify Routines

If, on the other hand, **notify** routines have been registered throughout the chain, then each interrupting device (by calling the **selnotify** kernel service) invokes the **notify** routine for the device above it in the calling chain. Thus in the preceding example, the **selnotify** kernel service for Device C calls the **notify** routine registered when Device B's **ddselect** routine invoked the **fp\_select** kernel service. Device B's **notify** routine must then decide whether to again call the **selnotify** kernel service to alert Device A's **notify** routine. If so, then Device A's **notify** routine is called, and makes its own determination whether to call another **selnotify** routine. If it does, the **selnotify** kernel service wakes up all the processes waiting on occurred events for Device A.

A variation on this scenario involves a cascaded chain in which only some device drivers have registered **notify** routines. In this case, the **selnotify** kernel service at each level calls the **notify** routine for the level above, until a level is encountered for which no **notify** routine was registered. At this point, all events of interest are determined to be true for the device driver at the head of the cascading chain. If any **notify** routines were registered in levels above the current level, they are never called.

### Returning from the fp\_select Kernel Service

The **fp\_select** kernel service does not wait for any selected events to become true, but returns immediately after the call to the object's **ddselect** routine has completed.

If the object's select routine is successfully called, the return code for the **fp\_select** kernel service is set to the return code provided by the object's **ddselect** routine.

### Execution Environment

The **fp\_select** kernel service can be called from the process environment only.

### Return Values

Item	Description
0	Indicates successful completion.
EAGAIN	Indicates that the allocation of internal data structures failed. The <i>rtneventp</i> parameter is not updated.
EINVAL	Indicates that the <i>fp</i> parameter is not a valid file pointer. The <i>rtneventp</i> parameter has the <b>POLLNVAL</b> flag set.

The **fp\_select** kernel service can also be set to the nonzero return code from the specified object's **ddselect** routine. The *rtneventp* parameter has the **POLLERR** flag set.

#### Related reference:

“fp\_poll Kernel Service” on page 161

“fp\_select Kernel Service notify Routine” on page 170

#### Related information:

select subroutine

## fp\_select Kernel Service notify Routine

### Purpose

Registers the **notify** routine.

### Syntax

```
#include <sys/types.h> #include <sys/errno.h> void notify (id, sub_id, rtnevents, pid) int id; int sub_id ;  
ushort rtnevents ; pid_t pid;
```

### Parameters

Item	Description
<i>id</i>	Indicates the selected function ID specified by the routine that made the call to the <b>selnotify</b> kernel service to indicate the occurrence of an outstanding event. For device drivers, this parameter is equivalent to the <i>devno</i> (device major and minor number) parameter.
<i>sub_id</i>	Indicates the unique ID specified by the routine that made the call to the <b>selnotify</b> kernel service to indicate the occurrence of an outstanding event. For device drivers, this parameter is equivalent to the <i>chan</i> parameter: channel for multiplexed drivers; 0 for nonmultiplexed drivers.
<i>rtnevents</i>	Specifies the <i>rtnevents</i> parameter supplied by the routine that made the call to the <b>selnotify</b> service indicating which events are designated as true.
<i>pid</i>	Specifies the process ID of a process waiting for the event corresponding to this call of the <b>notify</b> routine.

When a **notify** routine is provided for a cascaded function, the **selnotify** kernel service calls the specified **notify** routine instead of posting the process that was waiting on the event. It is up to this **notify** routine to determine if another **selnotify** call should be made to notify the waiting process of an event.

The **notify** routine is not called if the request is synchronous (that is, if the **POLLSYNC** flag is set in the *events* parameter) or if the original poll or select request is no longer outstanding.

**Note:** When more than one process has requested notification of an event and the **fp\_select** kernel service is used with a **notify** routine specified, the notification of the event causes the **notify** routine to be called once for each process that is currently waiting on one or more of the occurring events.

### Description

The **fp\_select** kernel service **notify** routine is registered by the caller of the **fp\_select** kernel service to be called by the kernel when specified events become true. The option to register this **notify** routine is available in a cascaded environment. The **notify** routine can be called at interrupt time.

### Execution Environment

The **fp\_select** kernel service **notify** routine can be called from either the process or interrupt environment.

#### Related reference:

“fp\_select Kernel Service” on page 167

“selnotify Kernel Service” on page 463

#### Related information:

Logical File System Kernel Services

## fp\_setea Kernel Service

### Purpose

Sets an extended attribute value.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
int fp_setea (fp, name, value, size, flags, segflag)
struct file *fp;
const char * name;
void * value;
size_t size;
int flags
int segflag;
```

## Parameters

Item	Description
<i>fp</i>	Specifies a file structure returned by the <b>fp_open</b> kernel service.
<i>name</i>	Specifies the name of the extended attribute. An extended attribute name is a NULL terminated string.
<i>value</i>	Specifies a pointer to the value of an attribute. The value of an extended attribute is an opaque byte stream of specified length.
<i>size</i>	Specifies the length of the value.
<i>flags</i>	None of the flags are defined at this time.
<i>segflag</i>	Specifies the flag indicating where the pointer specified by the <i>path</i> parameter is located: <b>SYS_ADSPACE</b> The pointers specified by the <i>name</i> and <i>value</i> parameters are stored in kernel memory. <b>USER_ADSPACE</b> The pointers specified by the <i>name</i> and <i>value</i> parameters are stored in application memory.

## Description

The **fp\_setea** kernel service provides a common service used by the following routines:

- The file system for the implementation of the **fsetea** subroutine
- Kernel routines outside the file system that set extended attribute values

## Execution Environment

The **fp\_setea** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
ERRNO	Indicates a failed operation. Returns an error number from the <code>/usr/include/sys/errno.h</code> file on failure.

### Related information:

fsetea subroutine

fp\_open subroutine

Logical File System Kernel Services

## fp\_statea Kernel Service

### Purpose

Provides information on an extended attribute.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
int fp_statea ( fp, name, buffer, segflag)
struct file * fp;
const char * name;
struct stat64x * buffer;
int segflag;
```

## Parameters

Item	Description
<i>fp</i>	Specifies a file structure returned by the <b>fp_open</b> kernel service.
<i>name</i>	Specifies the name of the extended attribute. An extended attribute name is a NULL terminated string.
<i>buffer</i>	Specifies a pointer to the <b>stat</b> structure in which information is returned.
<i>segflag</i>	Specifies the flag indicating the location of the pointer stored by the <i>path</i> parameter is located: <b>SYS_ADSPACE</b> The pointers specified by the <i>name</i> and <i>value</i> parameters are stored in kernel memory. <b>USER_ADSPACE</b> The pointers specified by the <i>name</i> and <i>value</i> parameters are stored in application memory.

## Description

The **fp\_statea** kernel service provides a common service used by the following routines:

- The file system for the implementation of the **fstatea** subroutine
- Kernel routines outside the file system that set extended attribute values.

## Execution Environment

The **fp\_statea** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
ERRNO	Indicates a failed operation. Returns an error number from the <code>/usr/include/sys/errno.h</code> file on failure.

### Related information:

fstatea subroutine

Logical File System Kernel Services

## fp\_write Kernel Service

### Purpose

Performs a write operation on an open file with arguments passed.

## Syntax

```
#include <sys/types.h> #include <sys/errno.h> int fp_write (fp, buf, nbytes, ext, seg, countp) struct file *
fp; char * buf; ssize_t nbytes, ext_t ext; int seg; ssize_t * countp;
```

## Parameters

Item	Description
<i>fp</i>	Points to a file structure returned by the <b>fp_open</b> or <b>fp_opendev</b> kernel service.
<i>buf</i>	Points to the buffer where data to be written to a file is located.
<i>nbytes</i>	Indicates the number of bytes to be written to the file.
<i>ext</i>	Specifies an extension argument required by some device drivers. Its content, form, and use are determined by the individual driver.
<i>seg</i>	Indicates in which part of memory the buffer specified by the <i>buf</i> parameter is located: <b>SYS_ADSpace</b> The buffer specified by the <i>buf</i> parameter is in kernel memory. <b>USER_ADSpace</b> The buffer specified by the <i>buf</i> parameter is in application memory.
<i>countp</i>	Points to the location where count of bytes actually written to the file is to be returned.

## Description

The **fp\_write** kernel service is an internal interface to the function provided by the **write** subroutine.

## Execution Environment

The **fp\_write** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
ERRNO	Returns an error number from the <code>/usr/include/sys/errno.h</code> file on failure.

### Related information:

write subroutine

Logical File System Kernel Services

## fp\_write Kernel Service for Data Link Control (DLC) Devices Purpose

Allows kernel data to be sent using a file pointer.

## Syntax

```
#include <sys/gdlexcb.h> #include <sys/fp_io.h> int fp_write (fp, buf, nbytes, ext, segflag, countp)
```

## Parameters

Item	Description
<i>fp</i>	Specifies file pointer returned from the <b>fp_open</b> kernel service.
<i>buf</i>	Points to a kernel <b>mbuf</b> structure.
<i>nbytes</i>	Contains the byte length of the write data. It is not necessary to set this field to the actual length of write data, however, since the <b>mbuf</b> contains a length field. Instead, this field can be set to any non-negative value (generally set to 0).
<i>ext</i>	Specifies the extended kernel service parameter. This is a pointer to the <b>dlc_io_ext</b> extended I/O structure for writes. The argument for this parameter must be in the kernel space. For more information on this parameter, see "write Subroutine Extended Parameters for DLC" <i>Technical Reference: Communications, Volume 1</i> .

<b>Item</b>	<b>Description</b>
<i>segflag</i>	Specifies the segment flag indicating where the <i>path</i> parameter is located. The only valid value is:
<i>countp</i>	<b>FP_SYS</b> The <i>path</i> parameter is stored in kernel memory. Points to the location where a count of bytes actually written is to be returned (must be in kernel space). GDLC does not provide this information for a kernel user since <b>mbufs</b> are used, but the file system requires a valid address and writes a copy of the <i>nbytes</i> parameter to that location.

## Description

Four types of data can be sent to generic data link control (GDLC). Network data can be sent to a service access point (SAP), and normal, exchange identification (XID) or datagram data can be sent to a link station (LS).

Kernel users pass a communications memory buffer (**mbuf**) directly to GDLC on the **fp\_write** kernel service. In this case, a **uiomove** kernel service is not required, and maximum performance can be achieved by merely passing the buffer pointer to GDLC. Each write buffer is required to have the proper buffer header information and enough space for the data link headers to be inserted. A write data offset is passed back to the kernel user at start LS completion for this purpose.

All data must fit into a single packet for each write call. That is, GDLC does not separate the user's write data area into multiple transmit packets. A maximum write data size is passed back to the user at **DLC\_ENABLE\_SAP** completion and at **DLC\_START\_LS** completion for this purpose.

Normally, a write subroutine can be satisfied immediately by GDLC by completing the data link headers and sending the transmit packet down to the device handler. In some cases, however, transmit packets can be blocked by the particular protocol's flow control or a resource outage. GDLC reacts to this differently, based on the system blocked/nonblocked file status flags (set by the file system and based on the **O\_NDELAY** and **O\_NONBLOCKED** values passed on the **fp\_open** kernel service). Nonblocked **write** subroutines that cannot get enough resources to queue the communications memory buffer (**mbuf**) return an error indication. Blocked write subroutines put the calling process to sleep until the resources free up or an error occurs. Each GDLC supports the **fp\_write** kernel service via its **dlcwrite** entry point. The **fp\_write** kernel service may be called from the process environment only.

## Return Values

<b>Item</b>	<b>Description</b>
0	Indicates a successful operation.
EAGAIN	Indicates that transmit is temporarily blocked, and the calling process cannot be put to sleep.

<b>Item</b>	<b>Description</b>
EINTR	Indicates that a signal interrupted the kernel service before it could complete successfully.
EINVAL	Indicates an invalid argument, such as too much data for a single packet.
ENXIO	Indicates an invalid file pointer.

These return values are defined in the **/usr/include/sys/errno.h** file.

### Related reference:

"fp\_write Kernel Service" on page 172

### Related information:

Generic Data Link Control (GDLC) Environment Overview

Parameter Blocks by ioctl Operation for DLC

## fp\_writew Kernel Service

### Purpose

Performs a write operation on an open file with arguments passed in **iovec** elements.

### Syntax

```
#include <sys/types.h> #include <sys/errno.h> int fp_writew (fp, iov, iovcnt, ext, seg, countp) struct file * fp; struct iovec * iov; ssize_t iovcnt; ext_t ext; int seg; ssize_t * countp;
```

### Parameters

Item	Description
<i>fp</i>	Points to a file structure returned by the <b>fp_open</b> kernel service.
<i>iov</i>	Points to an array of <b>iovec</b> elements. Each <b>iovec</b> element describes a buffer containing data to be written to the file.
<i>iovcnt</i>	Specifies the number of <b>iovec</b> elements in an array pointed to by the <i>iov</i> parameter.
<i>ext</i>	Specifies an extension argument required by some device drivers. Its content, form, and use are determined by the individual driver.
<i>segflag</i>	Indicates which part of memory the information designated by the <i>iov</i> parameter is located in: <b>SYS_ADSPACE</b> The information designated by the <i>iov</i> parameter is in kernel memory. <b>USER_ADSPACE</b> The information designated by the <i>iov</i> parameter is in application memory.
<i>countp</i>	Points to the location where the count of bytes actually written to the file is to be returned.

### Description

The **fp\_writew** kernel service is an internal interface to the function provided by the **writew** subroutine.

### Execution Environment

The **fp\_writew** kernel service can be called from the process environment only.

### Return Values

Item	Description
0	Indicates a successful operation.

If an error occurs, one of the values from the `/usr/include/sys/errno.h` file is returned.

### Related information:

writew subroutine

Logical File System Kernel Services

## fskv\_reg Kernel Service

### Purpose

Registers callout handlers for validation of file system operations.

### Syntax

```
#include <sys/xfops.h>  
  
int fskv_reg(fskv_t *fs_kv, ulong options);
```

```

typedef struct fskv {
    int version_number;
    int (*kv_open)(struct xfid *xfp,
                  long flags,
                  cred_ext_t *crxp);
    int (*kv_setattr)(struct xfid *xfp,
                     long op,
                     long arg1,
                     long arg2,
                     long arg3,
                     cred_ext_t *crxp);
} fskv_t;

```

## Parameters

### **fs\_kv**

Specifies an array of callout functions that are called to validate file system operations in the kernel.

### **options**

Specifies a bit mask of registration options. The **options** parameter is not defined currently. The caller must set the **options** parameter to 0.

## Description

The `fskv_reg` kernel service registers an array of functions that are called before the execution of file system-specific operations.

After a callout handler is registered, each of the affected operations is preceded by a call to the corresponding validation routine.

Only one callout array can be registered. After a callout array is registered with the `fskv_reg` kernel service, the subsequent invocation of the `fskv_reg` kernel service does not succeed until the `fskv_unreg` kernel service is called. The caller of the `fskv_reg` kernel service must have root authority.

## Execution environment

The `fskv_reg` kernel service can be called only from the process environment.

## Return values

On successful completion, the `fskv_reg` kernel service returns a value of 0.

The following error codes can be returned on failure:

### **EEXIST**

The callout array is already registered.

### **EPERM**

The caller does not have permission to invoke this function.

### **EINVAL**

A parameter is invalid.

## Callout handlers

Callouts can be specified for the `open`, `chmod`, `chown`, and `utimes` system calls. The `chmod`, `chown`, and `utimes` system calls are handled in a single operation in the kernel with the `setattr` call.

If the validation callout routine returns a nonzero value, the file system operation is stopped and the system call returns the `EPERM` value.



The validation routines are called only for local physical file systems (JFS2 and JFS) and network file system (NFS)-mounted file systems. The callout functions accept the following arguments. The `xfid` argument uniquely identifies the file within the current running system.

```
typedef struct xfid {
    fsid_t x_fsid;
    fid_t x_fid;
} xfid_t;
```

## **kv\_open() callout function**

The `kv_open` callout function contains the information that is available to the open routines of the file system to track and validate open calls.

### **Syntax**

```
#include <sys/file.h>
#include <sys/cred.h>

int (*kv_open)(struct xfid *xfp,
               long flags,
               void *nrp,
               cred_ext_t *crxp);
```

### **Parameters**

#### **xfp**

Pointer to an `xfid` structure that identifies the file system and object.

#### **nrp**

Name resolution information. If the `xfidToName` kernel service is called, this parameter must be passed to it. This pointer is not valid after it is returned from the callout function.

#### **flags**

Open flags that are passed by the application.

#### **crxp**

Pointer to the credentials for the calling process.

### **Return values**

#### **Zero**

Indicates that the validation completed successfully.

#### **Nonzero**

Indicates that the validation failed.

## **kv\_setattr() callout function**

The `kv_setattr` callout function contains the information that is available to the system call that initiated this function. The `setattr` function is called by the `chown`, `chmod`, and `utimes` system calls and the variants of those system calls (for example, `fchown`, `fchmod`, and `futimens` system calls).

### **Syntax**

```
#include <sys/vattr.h>
#include <sys/cred.h>

int (*kv_setattr)(struct xfid *xfp,
                  long op,
                  long arg1,
```

```

long arg2,
long arg3,
void *nrp,
cred_ext_t *crxp);

```

## Parameters

### xfp

Pointer to an xfid structure that identifies the file system and object.

**op** Specifies one of the following operations:

#### V\_OWN

Sets file ownership.

#### V\_MODE

Sets file mode.

#### V\_UTIME

Sets the file time specified by the user.

#### V\_STIME

Sets the file time requested by the system.

### argn

Specifies the following values for each of the listed operations.

Table 1. kv\_setattr() callout function: argn parameter values

Operations	arg1	arg2	arg3
V_OWN	flag: T_OWNER_AS_IS T_GROUP_AS_IS  (For information about the file ownership changes, see chownx subroutine.)	uid_t newuid	gid_t newgid
V_MODE	mode_t newmode	Unused	Unused
V_UTIME	flag: V_SETTIME  Ignore arguments and set time to the current time.	timestruct_t *atime  Set the access time.	timestruct_t *mtime  Set the modification time.
V_STIME	NULL or timestruct_t *atime  Set the access time.	NULL or timestruct_t *mtime  Set the modification time.	NULL or timestruct_t *ctime  Set the change time.

### nrp

Indicates the name resolution information. If the xfidToName kernel service is called, this parameter must be passed to it. This parameter is a pointer to temporary information that is not valid after it is returned from the validation routine.

### crxp

Pointer to credentials for the calling process.

## Return values

### Zero

Indicates that the validation completed successfully.

### Nonzero

Indicates that the validation failed.

**Related reference:**

“nameToXfid() Kernel Service” on page 377

“xfidToName() Kernel Service” on page 591

**fskv\_unreg Kernel Service****Purpose**

Unregisters callout handlers for validation of file system operations.

**Syntax**

```
#include <sys/xfops.h>
```

```
int fskv_unreg(ulong options);
```

**Parameters****options**

Specifies a bit mask of registration options. The **options** parameter is not defined currently. The caller must set the **options** parameter to 0.

**Description**

The fskv\_unreg kernel service removes the registration of the functions to be called before the execution of file system operations. After this service completes, the callout handlers are not executed.

The caller of the fskv\_unreg kernel service must have root authority.

**Execution environment**

The fskv\_unreg kernel service can be called only from the process environment.

**Return values**

On successful completion, the fskv\_unreg kernel service returns a value of 0.

The following error codes are returned on failure:

**EPERM**

The caller does not have permission to invoke this function.

**EINVAL**

A parameter is invalid.

**fubyte Kernel Service****Purpose**

Retrieves a byte of data from user memory.

**Syntax**

```
#include <sys/types.h> #include <sys/errno.h> int fubyte ( uaddr) uchar *uaddr;
```

**Parameter**

Item	Description
<i>uaddr</i>	Specifies the address of the user data.

## Description

The **fuword** kernel service fetches, or retrieves, a byte of data from the specified address in user memory. It is provided so that system calls and device heads can safely access user data. The **fuword** service ensures that the user has the appropriate authority to:

- Access the data.
- Protect the operating system from paging I/O errors on user data.

The **fuword** service should be called only while executing in kernel mode in the user process.

## Execution Environment

The **fuword** kernel service can be called from the process environment only.

## Return Values

When successful, the **fuword** service returns the specified byte.

Item	Description
-1	Indicates a <i>uaddr</i> parameter that is not valid.

The access is not valid under the following circumstances:

- The user does not have sufficient authority to access the data.
- The address is not valid.
- An I/O error occurs while referencing the user data.

### Related reference:

“fuword Kernel Service”

“subbyte Kernel Service” on page 478

### Related information:

Accessing User-Mode Data while in Kernel Mode

## fuword Kernel Service

### Purpose

Retrieves a word of data from user memory.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
int fuword ( uaddr)
int *uaddr;
```

### Parameter

Item	Description
<i>uaddr</i>	Specifies the address of user data.

## Description

The **fuword** kernel service retrieves a word of data from the specified address in user memory. It is provided so that system calls and device heads can safely access user data. The **fuword** service ensures that the user had the appropriate authority to:

- Access the data.
- Protect the operating system from paging I/O errors on user data.

The **fuword** service should be called only while executing in kernel mode in the user process.

## Execution Environment

The **fuword** kernel service can be called from the process environment only.

## Return Values

When successful, the **fuword** service returns the specified word of data.

Item	Description
-1	Indicates a <i>uaddr</i> parameter that is not valid.

The access is not valid under the following circumstances:

- The user does not have sufficient authority to access the data.
- The address is not valid.
- An I/O error occurred while referencing the user data.

For the **fuword** service, a retrieved value of -1 and a return code of -1 are indistinguishable.

### Related reference:

“fubyte Kernel Service” on page 179

“subyte Kernel Service” on page 478

### Related information:

Accessing User-Mode Data while in Kernel Mode

## g

The following kernel services begin with the with the letter g.

### getblk Kernel Service

#### Purpose

Assigns a buffer to the specified block.

#### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/buf.h>

struct buf *getblk
( dev, blkno)
dev_t dev;
daddr_t blkno;
```

## Parameters

Item	Description
<i>dev</i>	Specifies the device that contains the block to be allocated.
<i>blkno</i>	Specifies the block to be allocated.

## Description

The **getblk** kernel service first checks whether the specified buffer is in the buffer cache. If the buffer resides there, but is in use, the **e\_sleep** service is called to wait until the buffer is no longer in use. Upon waking, the **getblk** service tries again to access the buffer. If the buffer is in the cache and not in use, it is removed from the free list and marked as busy. Its buffer header is then returned. If the buffer is not in the buffer cache, another buffer is taken from the free list and returned.

## Execution Environment

The **getblk** kernel service can be called from the process environment only.

## Return Values

The **getblk** service returns a pointer to the buffer header. A nonzero value for **B\_ERROR** in the **b\_flags** field of the buffer header (**bufstructure**) indicates an error. If this occurs, the caller should release the block's buffer using the **brelease** kernel service.

On a platform that supports storage keys, the buffer header is allocated from the storage that is protected by the **KKEY\_BLOCK\_DEV** kernel key.

### Related reference:

“bread Kernel Service” on page 30

### Related information:

Block I/O Buffer Cache Kernel Services: Overview

I/O Kernel Services

## getc Kernel Service

### Purpose

Retrieves a character from a character list.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <cblock.h>
```

```
int getc ( header)
struct clist *header;
```

### Parameter

Item	Description
<i>header</i>	Specifies the address of the <b>clist</b> structure that describes the character list.

## Description

**Attention:** The caller of the **getc** service must ensure that the character list is pinned. This includes the **clist** header and all the **cblock** character buffers. Otherwise, the system may crash.

The **getc** kernel service returns the character at the front of the character list. After returning the last character in the buffer, the **getc** service frees that buffer.

## Execution Environment

The **getc** kernel service can be called from either the process or interrupt environment.

## Return Values

Item	Description
-1	Indicates that the character list is empty.

## Related information:

I/O Kernel Services

## **getc** Kernel Service Purpose

Removes the first buffer from a character list and returns the address of the removed buffer.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <cblock.h>
```

```
struct cblock *getc
( header)
struct clist *header;
```

## Parameter

Item	Description
<i>header</i>	Specifies the address of the <b>clist</b> structure that describes the character list.

## Description

**Attention:** The caller of the **getc** service must ensure that the character list is pinned. This includes the **clist** header and all the **cblock** character buffers. Character buffers acquired from the **getc** service are pinned. Otherwise, the system may crash.

The **getc** kernel service returns the address of the character buffer at the start of the character list and removes that buffer from the character list. The user must free the buffer with the **putc** service when finished with it.

## Execution Environment

The **getc** kernel service can be called from either the process or interrupt environment.

## Return Values

A null address indicates the character list is empty.

The **getcb** service returns the address of the character buffer at the start of the character list when the character list is not empty.

### Related reference:

“getcf Kernel Service” on page 185

### Related information:

I/O Kernel Services

## getcbp Kernel Service

### Purpose

Retrieves multiple characters from a character buffer and places them at a designated address.

### Syntax

```
#include <cblock.h>
```

```
int getcbp ( header, dest, n)
struct clist *header;
char *dest;
int n;
```

### Parameters

Item	Description
<i>header</i>	Specifies the address of the <b>clist</b> structure that describes the character list.
<i>dest</i>	Specifies the address where the characters obtained from the character list are to be placed.
<i>n</i>	Specifies the number of characters to be read from the character list.

### Description

**Attention:** The caller of the **getcbp** services must ensure that the character list is pinned. This includes the **clist** header and all the **cblock** character buffers. Character buffers acquired from the **getcf** service are pinned. Otherwise, the system may crash.

The **getcbp** kernel service retrieves as many as possible of the *n* characters requested from the character buffer at the start of the character list. The **getcbp** service then places them at the address pointed to by the *dest* parameter.

### Execution Environment

The **getcbp** kernel service can be called from either the process or interrupt environment.

### Return Values

The **getcbp** service returns the number of characters retrieved from the character buffer.

### Related reference:

“getcf Kernel Service” on page 185

### Related information:

I/O Kernel Services



## **getc Kernel Service**

### **Purpose**

Retrieves a free character buffer.

### **Syntax**

```
#include <sys/types.h>
#include <sys/errno.h>
#include <cblock.h>
struct cblock *getc ( )
```

### **Description**

The **getc** kernel service retrieves a character buffer from the list of available ones and returns that buffer's address. The returned character buffer is pinned. If you use the **getc** service to get a character buffer, be sure to free the space when you have finished using it. The buffers received from the **getc** service should be freed by using the **putc** kernel service.

Before starting the **getc** service, the caller should request enough **clist** resources by using the **pincf** kernel service. The proper use of the **getc** service ensures that there are sufficient pinned buffers available to the caller.

If the **getc** service indicates that there is no available character buffer, the **waitcfree** service can be called to wait until a character buffer becomes available.

The **getc** service has no parameters.

### **Execution Environment**

The **getc** kernel service can be called from either the process or interrupt environment.

### **Return Values**

Upon successful completion, the **getc** service returns the address of the allocated character buffer.

A null pointer indicates no buffers are available.

#### **Related reference:**

“pincf Kernel Service” on page 410

“putc Kernel Service” on page 426

#### **Related information:**

I/O Kernel Services

## **getc Kernel Service**

### **Purpose**

Returns the character at the end of a designated list.

### **Syntax**

```
#include <sys/types.h>
#include <sys/errno.h>
#include <cblock.h>
```

```
int getc ( header)
struct clist *header;
```

## Parameter

Item	Description
<i>header</i>	Specifies the address of the <b>clist</b> structure that describes the character list.

## Description

**Attention:** The caller of the **getc** service must ensure that the character list is pinned. This includes the **clist** header and all the **cblock** character buffers. Character buffers acquired from the **getc** service are pinned.

The **getc** kernel service is identical to the **getc** service, except that the **getc** service returns the character at the end of the list instead of the character at the front of the list. The character at the end of the list is the last character in the first buffer, not in the last buffer.

## Execution Environment

The **getc** kernel service can be called from either the process or interrupt environment.

## Return Values

The **getc** service returns the character at the end of the list instead of the character at the front of the list.

### Related reference:

“getc Kernel Service” on page 185

### Related information:

I/O Kernel Services

## getblk Kernel Service Purpose

Allocates a free buffer.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/buf.h>
struct buf *getblk ( )
```

## Description

**Attention:** The use of the **getblk** service by character device drivers is strongly discouraged. As an alternative, character device drivers can use the **xmalloc** service to allocate the memory space directly, or the character I/O kernel services such as the **getc** or **getc** services.

The **getblk** kernel service allocates a buffer and buffer header and returns the address of the buffer header. If no free buffers are available, then the **getblk** service waits for one to become available. Block device drivers can retrieve buffers using the **getblk** service.

In the header, the **b\_forw**, **b\_back**, **b\_flags**, **b\_bcount**, **b\_dev**, and **b\_un** fields are used by the system and cannot be modified by the driver. The **av\_forw** and **av\_back** fields are available to the user of the **getblk** service for keeping a chain of buffers by the user of the **getblk** service. (This user could be the kernel file system or a device driver.) The **b\_blkno** and **b\_resid** fields can be used for any purpose.

The **brfree** service is used to free this type of buffer.

The `geteblk` service has no parameters.

## Execution Environment

The `geteblk` kernel service can be called from the process environment only.

## Return Values

The `geteblk` service returns a pointer to the buffer header. There are no error codes because the `geteblk` service waits until a buffer header becomes available.

On a platform that supports storage keys, the buffer header is allocated from the storage protected by the `KKEY_BLOCK_DEV` kernel key.

### Related reference:

“`xmalloc` Kernel Service” on page 598

“`buf` Structure” on page 615

### Related information:

Block I/O Buffer Cache Kernel Services: Overview

## geterror Kernel Service Purpose

Determines the completion status of the buffer.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/buf.h>
```

```
int geterror ( bp)
struct buf *bp;
```

## Parameter

Item	Description
<i>bp</i>	Specifies the address of the buffer structure whose status is to be checked.

On a platform that supports storage keys, the passed in *bp* parameter must be in the `KKEY_PUBLIC` or `KKEY_BLOCK_DEV` protection domain.

## Description

The `geterror` kernel service checks the specified buffer to see if the `b_error` flag is set. If that flag is not set, the `geterror` service returns 0. Otherwise, it returns the nonzero `B_ERROR` value or the `EIO` value (if `b_error` is 0).

## Execution Environment

The `geterror` kernel service can be called from either the process or interrupt environment.

## Return Values

Item	Description
0	Indicates that no I/O error occurred on the buffer.
b_error value	Indicates that an I/O error occurred on the buffer.
EIO	Indicates that an unknown I/O error occurred on the buffer.

#### Related information:

Block I/O Buffer Cache Kernel Services: Overview

I/O Kernel Services

## getexcept Kernel Service

### Purpose

Allows kernel exception handlers to retrieve additional exception information.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/except.h>
```

```
void getexcept
( exceptp)
struct except *exceptp;
```

### Parameter

Item	Description
<i>exceptp</i>	Specifies the address of an <b>except</b> structure, as defined in the <code>/usr/include/sys/except.h</code> file. The <b>getexcept</b> service copies detailed exception data from the current machine-state save area into this caller-supplied structure.

### Description

The **getexcept** kernel service provides exception handlers the capability to retrieve additional information concerning the exception from the machine-state save area.

The **getexcept** service should only be used by exception handlers when called to handle an exception. The contents of the structure pointed at by the *exceptp* parameter is platform-specific, but is described in the `/usr/include/sys/except.h` file for each type of exception that provides additional data. This data is typically included in any error logging data for the exception. It can be also used to attempt to handle or recover from the exception.

### Execution Environment

The **getexcept** kernel service can be called from either the process or interrupt environment. It should be called only when handling an exception.

### Return Values

The **getexcept** service has no return values.

#### Related information:

Kernel Extension and Device Driver Management Kernel Services and

## getfslimit Kernel Service

### Purpose

Returns the maximum file size limit of the current process.

## Syntax

```
#include <sys/types.h> offset_t getfslimit (void)
```

## Description

The **getfslimit** kernel service returns the file size limit of the current process as a 64 bit integer. This can be used by file systems to implement the checks needed to enforce limits. The **getfslimit** kernel service is called from the process environment.

## Return Values

The **getfslimit** kernel service returns the the file size limit, there are no error values.

### Related information:

ulimit subroutine

getrlimit subroutine

ulimit subroutine

## get\_pag or get\_pag64 Kernel Service Purpose

Retrieves a Process Authentication Group (PAG) value for the current process.

## Syntax

```
#include <sys/cred.h>
```

```
int get_pag ( type, pag )
int type;
int *pag;
```

```
int get_pag64 ( type, pag )
int type;
uint64_t *pag;
```

## Parameters

Item	Description
<i>type</i>	PAG type to retrieve
<i>pag</i>	Pointer to buffer where operating system returns the PAG

## Description

The **get\_pag** and **get\_pag64** kernel services copy the requested PAG from the current process into *pag*. The value of *type* must be a defined PAG ID. The PAG ID for the Distributed Computing Environment (DCE) is 0.

## Execution Environment

The **get\_pag** and **get\_pag64** kernel services can be called from the process environment only.

## Return Values

A value of 0 is returned upon successful completion. If unsuccessful, **errno** is set to a value that explains the error.

## Error Codes

The `get_pag` kernel service fails if one or both of the following conditions are true:

Item	Description
EINVAL	Invalid PAG specification
E_OVERFLOW	PAG value is 64-bit (should be using <code>get_pag64</code> )

The `get_pag64` kernel service fails if the following condition is true:

Item	Description
EINVAL	Invalid PAG specification

### Related information:

Security Kernel Services

## getpid Kernel Service

### Purpose

Gets the process ID of the current process.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
pid_t getpid ()
```

### Description

The `getpid` kernel service returns the process ID of the calling process.

The `getpid` service can also be used to check the environment that the routine is being executed in. If the caller is executing in the interrupt environment, the `getpid` service returns a process ID of -1. If a routine is executing in a process environment, the `getpid` service obtains the current process ID.

### Execution Environment

The `getpid` kernel service can be called from either the process or interrupt environment.

### Return Values

Item	Description
-1	Indicates that the <code>getpid</code> service was called from an interrupt environment.

The `getpid` service returns the process ID of the current process if called from a process environment.

### Related information:

Process and Exception Management Kernel Services

Understanding Execution Environments

## getppid Kernel Service

### Purpose

Gets the parent process ID of the specified process.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
pid_t getppid (ProcessID)
pid_t ProcessID;
```

## Parameter

Item	Description
<b>ProcessID</b>	Specifies the process ID. If this parameter is 0, then the parent process ID of the calling process is returned.

## Description

The **getppid()** kernel service accepts a process ID as an input. If the input process ID is 0, the **getppid()** subroutine returns the process ID of the calling process' parent process. If the input process ID is nonzero and a valid value, the parent ID of the input process ID is returned. If the input process ID is invalid, the **getppid()** kernel service returns -1.

## Execution Environment

The **getppid()** kernel service can be called from the process environment only.

## Return Values

Item	Description
-1	Indicates that the <b>ProcessID</b> parameter is invalid.

### Related reference:

“getpid Kernel Service” on page 190

### Related information:

Process and Exception Management Kernel Services

Understanding Execution Environments

## getuerror Kernel Service Purpose

Allows kernel extensions to read the **ut\_error** field for the current thread.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
int getuerror ()
```

## Description

The **getuerror** kernel service allows a kernel extension in a process environment to retrieve the current value of the current thread's **ut\_error** field. Kernel extensions can use the **getuerror** service when using system calls or other kernel services that return error information in the **ut\_error** field.

For system calls, the system call handler copies the value of the **ut\_error** field in the per thread **uthread** structure to the **errno** global variable before returning to the caller. However, when kernel services use available system calls, the system call handler is bypassed. The **getuerror** service must then be used to obtain error information.

## Execution Environment

The **getuerror** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.

When an error occurs, the **getuerror** kernel service returns the current value of the **ut\_error** field in the per thread **uthread** structure. Possible return values for this field are defined in the **/usr/include/sys/errno.h** file.

### Related reference:

“setuerror Kernel Service” on page 469

### Related information:

Kernel Extension and Device Driver Management Kernel Services  
Understanding System Call Execution

## getufdflags and setufdflags Kernel Services

### Purpose

Queries and sets file-descriptor flags.

### Syntax

```
#include <sys/user.h> int getufdflags(fd, flagsp) int fd; int *flagsp; #include <sys/user.h> int setufdflags(fd, flags) int fd; int flags;
```

### Parameters

Item	Description
<i>fd</i>	Identifies the file descriptor.
<i>flags</i>	Sets attribute flags for the specified file descriptor. Refer to the <b>sys/user.h</b> file for the list of valid flags.
<i>flagsp</i>	Points to an integer field where the flags associated with the file descriptor are stored on successful return.

### Description

The **setufdflags** and **getufdflags** kernel services set and query the file descriptor flags. The file descriptor flags are listed in **fontl.h**.

## Execution Environment

These kernel services can be called from the process environment only.

## Return Values



Item	Description
0	Indicates successful completion.
EBADF	Indicates that the <i>fd</i> parameter is not a file descriptor for an open file.

#### Related reference:

“*ufdhold* and *ufdrele* Kernel Service” on page 516

## get\_umask Kernel Service

### Purpose

Queries the file mode creation mask.

### Syntax

```
int get_umask(void)
```

### Description

The *get\_umask* service gets the value of the file mode creation mask currently set for the process.

**Note:** There is no corresponding kernel service to set the umask because kernel routines that need to set the umask can call the *umask* subroutine.

### Execution Environment

The *get\_umask* kernel service can be called from the process environment only.

### Return Values

The *get\_umask* kernel service always completes successfully. Its return value is the current value of the *umask*.

#### Related information:

*umask* subroutine

## gfsadd Kernel Service

### Purpose

Adds a file system type to the *gfs* table.

### Syntax

```
#include <sys/types.h> #include <sys/errno.h> int gfsadd ( gfsno, gfsp) int gfsno; struct gfs *gfsp;
```

### Parameters

Item	Description
<i>gfsno</i>	Specifies the file system number. This small integer value is either defined in the <i>/usr/include/sys/vmount.h</i> file or a user-defined number of the same order.
<i>gfsp</i>	Points to the file system description structure.

### Description

The *gfsadd* kernel service is used during configuration of a file system. The configuration routine for a file system invokes the *gfsadd* kernel service with a *gfs* structure. This structure describes the file system type.

The **gfs** structure type is defined in the **/usr/include/sys/gfs.h** file. The **gfs** structure must have the following fields filled in:

<b>Field</b>	<b>Description</b>
<b>gfs_type</b>	Specifies the integer type value. The predefined types are listed in the <b>/usr/include/sys/vmount.h</b> file.
<b>gfs_name</b>	Specifies the character string name of the file system. The maximum length of this field is 16 bytes. Shorter names must be null-padded.
<b>gfs_flags</b>	Specifies the flags that define the capabilities of the file system. The following flag values are defined: <b>GFS_AHAFS_INFO</b> GFS supports AHAFS FS monitoring. <b>GFS_AIX_FLOCK</b> Uses <b>common_relock()</b> to manage advisory locks. <b>GFS_DIROP</b> Call parent <b>vnop</b> instead of <b>obj</b> . <b>GFS_FASTPATH</b> GFS supports AIO fast path. <b>GFS_FUMNT</b> File system supports forced unmount. <b>GFS_INIT</b> GFS has been initialized <b>GFS_MEMCNTL</b> New <b>memcntl</b> vnode operation <b>GFS_MLS</b> GFS supports MLS. <b>GFS_NAMED_OPEN</b> File system supports named open. <b>GFS_NO_ACCT</b> Do not do file system account on this file system. <b>GFS_NOEXPORT</b> GFS cannot be exported by NFS. <b>GFS_NOUMASK</b> File system does not apply umask when creating new objects. <b>GFS_OFLAGS64</b> GFS supports 64 bit open flags. <b>GFS_REMNT</b> File system supports remount of a mounted file system.

Field	Description
	<b>GFS_REMOTE</b> File system is remote (ie. NFS).
	<b>GFS_STATFSVP</b> File system supports <b>vfs_statfsvp</b> VFS interface. (new vfs operation: <b>vfs_statfsvp</b> )
	<b>GFS_SYS5DIR</b> File system that uses the System V-type directory structure.
	<b>GFS_SYNCVFS</b> The <b>syncvfs</b> vnode operation.
	<b>GFS_VERSION4</b> File system supports AIX Version 4 V-node interface.
	<b>GFS_VERSION42</b> File system supports AIX 4.2 V-node interface. (new vnode operation: <b>vnop_seek</b> )
	<b>GFS_VERSION421</b> File system supports AIX 4.2.1 V-node interface.(new vnode operations: <b>vnop_sync_range</b> , <b>vnop_create_attr</b> , <b>vnop_finfo</b> , <b>vnop_map_lloff</b> , <b>vnop_readdir_eofp</b> , <b>vnop_rdwr_attr</b> )
	<b>GFS_VERSION43</b> File system supports AIX 4.3 V-node interface. (new file flag for <b>vnop_sync_range</b> :FMSYNC)
	<b>GFS_VERSION53</b> File system supports AIX 5.3 V-node interface (new vnode operations: <b>vnop_getxcl</b> , <b>vnop_setxcl</b> ) and AIX 5.3 VFS interface. (new vfs operation: <b>vfs_aclxcntl</b> )
	<b>GFS_VREGSEL</b> GFS wants to select vnode operation called for <b>VREG</b> files.
<code>gfs_ops</code>	Specifies the array of pointers to <b>vfs</b> operation implementations.
<code>gn_ops</code>	Specifies the array of pointers to v-node operation implementations.

The file system description structure can also specify:

Item	Description
<code>gfs_init</code>	Points to an initialization routine to be called by the <b>gfsadd</b> kernel service. This field must be null if no initialization routine is to be called.
<code>gfs_data</code>	Points to file system private data.

## Execution Environment

The **gfsadd** kernel service can be called from the process environment only.

## Return Values

Item	Description
<code>0</code>	Indicates successful completion.
<code>EBUSY</code>	Indicates that the file system type has already been installed.
<code>EINVAL</code>	Indicates that the <i>gfsno</i> value is larger than the system-defined maximum. The system-defined maximum is indicated in the <code>/usr/include/sys/vmount.h</code> file.

### Related reference:

“**gfsdel** Kernel Service”

“**vfs\_init** Entry Point” on page 645

## **gfsdel** Kernel Service

### Purpose

Removes a file system type from the **gfs** table.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
int gfsdel ( gfsno)
int gfsno;
```

## Parameter

Item	Description
<i>gfsno</i>	Specifies the file system number. This value identifies the type of the file system to be deleted.

## Description

The **gfsdel** kernel service is called to delete a file system type. It is not valid to mount any file system of the given type after that type has been deleted.

## Execution Environment

The **gfsdel** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates successful completion.
ENOENT	Indicates that the indicated file system type was not installed.
EINVAL	Indicates that the <i>gfsno</i> value is larger than the system-defined maximum. The system-defined maximum is indicated in the <i>/usr/include/sys/vmount.h</i> file.
EBUSY	Indicates that there are active <b>vfs</b> structures for the file system type being deleted.

### Related reference:

“gfsadd Kernel Service” on page 193

### Related information:

Virtual File System Overview

Virtual File System Kernel Services

## gn\_closecnt Subroutine

### Purpose

Maintains the using count on a **gnode** structure.

## Syntax

```
#include <sys/vnode.h>
#include <sys/fcntl.h>
```

```
void gn_closecnt (gnode, flags)
struct gnode *gnode;
long flags;
```

## Parameters

Item	Description
<i>gnode</i>	Points to a <b>gnode</b> structure.
<i>flags</i>	Specifies the open mode ( <b>FREAD</b> , <b>FWRITE</b> , <b>FEXEC</b> , <b>FRSHARE</b> ) from the open file flags.

## Description

The **gn\_closecnt** subroutine uses the passed in *flags* value to determine the appropriate using counts to decrease in the *gnode* structure. For example, if the **FREAD** flag is set, the **gn\_closecnt** subroutine decreases the *gn\_rdcnt* field. The following table shows the mapping of the *flags* value to the counts field in the *gnode* structure:

Item	Description
<b>FREAD</b>	<i>gn_rdcnt</i>
<b>FWRITE</b>	<i>gn_wrcnt</i>
<b>FEXEC</b>	<i>gn_excnt</i>
<b>FRSHARE</b>	<i>gn_rshcnt</i>

## Return Values

The **gn\_closecnt** subroutine returns no return values.

## Error Codes

The **gn\_closecnt** subroutine returns no error codes.

### Related information:

Understanding Data Structures and Header Files for Virtual File Systems

## gn\_common\_memcntl Subroutine

### Purpose

Changes or queries the physical attachment of a file.

### Syntax

```
#include <sys/vnode.h>
#include <sys/fcntl.h>
```

```
int gn_common_memcntl (gnode, cmd, arg)
struct gnode * gnode;
int cmd;
void * arg;
```

### Parameters

Item	Description
<i>gnode</i>	Points to a <b>gnode</b> structure.
<i>cmd</i>	Specifies the operation to be performed. The <i>cmd</i> parameter can be one of the following values: <ul style="list-style-type: none"> <li>• <b>F_ATTACH</b></li> <li>• <b>F_DETACH</b></li> <li>• <b>F_ATTINFO</b></li> </ul>
<i>arg</i>	Points to a structure containing information for the specified <i>cmd</i> parameter. <ul style="list-style-type: none"> <li><b>F_ATTACH</b> attach_desc_t</li> <li><b>F_DETACH</b> detach_desc_t</li> <li><b>F_ATTINFO</b> attinfo_desc_t</li> </ul>

## Description

The `gn_common_memcntl` subroutine is to be called by file system `vnode_memcntl` implementations. It performs the normal function of such operations. If the `cmd` parameter is set to `F_ATTACH`, the `gn_common_memcntl` subroutine attaches the segment specified by the `gn_seg` field in the `gnode` structure. If the `cmd` parameter is set to `F_DETACH`, the `gn_common_memcntl` subroutine detaches the segment. If the `cmd` parameter is set to `F_ATTINFO`, the `gn_common_memcntl` subroutine returns information about the current state of attachment.

## Return Values

Item	Description
0	Success.
non-zero	Failure.

## Error Codes

Item	Description
EINVAL	The <code>cmd</code> parameter is not valid.
ENOMEM	Resources are not available to attach the memory segment.

## gn\_mapcnt Subroutine Purpose

Maintains the mapping count in a `gnode` structure.

## Syntax

```
#include <sys/vnode.h>
#include <sys/shm.h>
```

```
void gn_mapcnt (gnode, flags)
struct gnode * gnode;
long flags;
```

## Parameters

Item	Description
<i>gnode</i>	Points to a <code>gnode</code> structure.
<i>flags</i>	Specifies the following mapping flag:  <code>SHM_RDONLY</code> Only read access is required.

## Description

The `gn_mapcnt` subroutine uses the passed in `flags` value to determine the appropriate mapping count to increase in the `gnode` structure. If the `SHM_RDONLY` flag is set, the `gn_mapcnt` subroutine increases the `gn_mrdcnt` field. Otherwise, the `gn_mapcnt` subroutine increases the `gn_mwrcnt` field.

## Return Values

The `gn_mapcnt` subroutine returns no return values.

## Error Codes

The `gn_mapcnt` subroutine returns no error codes.

### Related information:

`mmap` subroutine

`shmat` subroutine

## gn\_opencnt Subroutine

### Purpose

Maintains the using count on a `gnode` structure.

### Syntax

```
#include <sys/vnode.h>
```

```
#include <sys/fcntl.h>
```

```
void gn_opencnt (gnode, flags)
```

```
struct gnode * gnode;
```

```
long flags;
```

### Parameters

Item	Description
<i>gnode</i>	Points to a <code>gnode</code> structure.
<i>flags</i>	Specifies the open mode ( <code>FREAD</code> , <code>FWRITE</code> , <code>FEXEC</code> , <code>FRSHARE</code> ) from the open file flags.

### Description

The `gn_opencnt` subroutine uses the passed in *flags* value to determine the appropriate using counts to increase in the *gnode* structure. The following table shows the mapping of the *flags* value to the counts field in the *gnode* structure:

Item	Description
<code>FREAD</code>	<code>gn_rdcnt</code>
<code>FWRITE</code>	<code>gn_wrcnt</code>
<code>FEXEC</code>	<code>gn_excnt</code>
<code>FRSHARE</code>	<code>gn_rshcnt</code>

### Return Values

The `gn_opencnt` subroutine returns no return values.

## Error Codes

The `gn_opencnt` subroutine returns no error codes.

### Related information:

Understanding Data Structures and Header Files for Virtual File Systems

## gn\_unmapcnt Subroutine

### Purpose

Maintains the mapping count in a `gnode` structure.

## Syntax

```
#include <sys/vnode.h>
#include <sys/shm.h>
```

```
void gn_unmapcnt (gnode, flags)
struct gnode * gnode;
long flags;
```

## Parameters

Item	Description
<i>gnode</i>	Points to a <b>gnode</b> structure.
<i>flags</i>	Specifies the following mapping flag: <b>SHM_RDONLY</b> Only read access is required.

## Description

The **gn\_unmapcnt** subroutine uses the passed in *flags* value to determine the appropriate mapping count to decrease in the *gnode* structure. If the **SHM\_RDONLY** flag is set, the **gn\_unmapcnt** subroutine decreases the *gn\_mrdcnt* field. Otherwise, the **gn\_unmapcnt** subroutine decreases the *gn\_mwrcnt* field.

## Return Values

The **gn\_unmapcnt** subroutine returns no return values.

## Error Codes

The **gn\_unmapcnt** subroutine returns no error codes.

### Related information:

mmap subroutine  
shmat subroutine

## groupmember, groupmember\_cr Subroutines Purpose

Determines if the named group is a member of a credential group set.

## Syntax

```
#include <sys/types.h>
#include <sys/cred.h>
```

```
int groupmember (gid)
gid_t gid;
```

```
int groupmember_cr (gid, cred)
gid_t gid;
cred_t * cred;
```

## Parameters



Item	Description
<i>gid</i>	Specifies an identifier for a group.
<i>cred</i>	Points to a <b>ucred</b> structure.

## Description

The **groupmember** subroutines determine if a group is included in the group set of a credential structure. The **groupmember** subroutine queries the credential associated with the current thread. The **groupmember\_cr** subroutine checks for the group within the specified **ucred** structure.

## Return Values

The **groupmember** subroutines return TRUE if the **ucred** structure contains the specified *gid* parameter or if the specified *gid* parameter is the current effective group ID for the thread. Otherwise, these routines return FALSE.

## Error Codes

The **groupmember** subroutines return no error codes.

### Related information:

Security Kernel Services

## h

The following kernel services begin with the with the letter h.

### heap\_create Kernel Service

#### Purpose

Initializes a new heap to be used with kernel memory management services. The **heap\_create** kernel service replaces the **init\_heap** kernel service. It returns a heap handle that can be used with the **xmalloc** and the **xmfree** kernel services.

#### Syntax

```
#include <sys/types.h>
#include <sys/malloc.h>
#include <sys/skeys.h>
#include <sys/kernno.h>
```

```
kernno_t heap_create (heapattr_t * heapattr, heapaddr_t * heapptr);
```

#### Parameters

Item	Description
<i>heapattr</i>	Points to an initialized heap attribute structure. See the <b>sys/malloc.h</b> file. This structure is initialized by the caller of <b>heap_create</b> .
<i>heapptr</i>	Points to an external heap descriptor. The caller must initialize this parameter to the <b>HPA_INVALID_HEAP</b> value.

The *heapattr* structure contains the following fields:

<b>Item</b>	<b>Description</b>
<i>eye_catch8b_t hpa_eyec</i>	Must be initialized to the <b>EYEC_HEAPATTR</b> value.
<i>short hpa_version</i>	Must be initialized to the <b>HPA_VERSION</b> value.
<i>long hpa_flags</i>	The following flags describe heap properties:
	<b>HPA_PAGED</b> The heap returns pageable memory.
	<b>HPA_PINNED</b> The heap returns pinned memory.
	<b>HPA_SHARED</b> The returned descriptor is backed by a common sub-heap.
	<b>HPA_PRIVATE</b> The returned descriptor is backed by isolated storage.
<i>void * hpa_heapaddr</i>	Must be set to NULL (reserved).
<i>size_t hpa_heapsize</i>	Heap size in bytes. It is only used for private heaps.
<i>size_t hpa_limit</i>	Usage barrier independent from size. Limits the amount available from a private heap that is less than or equal to the actual size of the private heap.
<i>long hpa_debug_level</i>	Heap debug level. The <b>HPA_DEFAULT_DEBUG</b> value gives the heap the system debug level.
<i>uint hpa_kkey</i>	Kernel key requested for the storage allocated.

## Description

The **heap\_create** service is a replacement for the **init\_heap** service. It can be used to create private heaps, and to create shared sub-heaps. After this service creates a private heap or a handle to a shared sub-heap, the returned **heapaddr\_t** value can be used with the **xmalloc** service or the **xmfree** service to allocate or free memory from that heap.

The most common usage for the **heap\_create** service is to get a handle to a shared sub-heap. This is done by setting the **HPA\_SHARED** flag in the input attribute structure. See the **sys\_malloc.h** file.

Private heaps can be created by specifying the **HPA\_PRIVATE** flag. This allows the **heap\_create** service to initialize and manage an area of virtual memory as a private heap. The **hpa\_heapaddr** field must be set to zero. The **heap\_create** service provides the storage but this field is reserved for future use. The **hpa\_size** field indicates the size of the private heap in bytes.

Private heaps can make use of the **hpa\_limit** field. Use the **hpa\_size** field to reserve a maximum effective address space. Then use the **hpa\_limit** field to alter and control the amount of effective address space that is in use. The value of the **hpa\_limit** field must be less than or equal to the value of the **hpa\_size** field.

The **hpa\_debug** and **hpa\_kkey** fields are required for shared and private heaps. The **hpa\_debug** level allows a component run-time debug level to be applied to allocations using the returned heap handle. The **hpa\_kkey** field associates a kernel key with a sub-heap that can limit the kernel accessibility.

On a successful completion, the *heapattr* field contains the address of a heap structure. This can be used as a parameter to the **xmalloc** and the **xmfree** kernel services. The memory returned by these services and the internal heap structures are protected by the **hpa\_kkey** field. When calling the **xmalloc** and the **xmfree** heap services, the caller must hold the key that was used when creating the heap.

## Execution Environment

The **heap\_create** kernel service can be called from the process environment only.

## Return Values

Item	Description
0 EINVAL_HEAP_CREATE	Indicates a successful completion. A descriptor is returned in the <i>heapptr</i> parameter. Indicates one or more of the following inputs that were not valid: <ul style="list-style-type: none"> <li>• <i>heapattr</i> is NULL.</li> <li>• <i>*heapptr</i> != HPA_INVALID_HEAP.</li> <li>• <i>heapattr-&gt;hpa_eyec</i> != EYEC_HEAPATTR.</li> <li>• <i>heapattr-&gt;hpa_version</i> != HPA_VERSION.</li> <li>• Flags: Both the HPA_SHARED and the HPA_PRIVATE flags are specified.</li> <li>• Flags: Neither the HPA_SHARED nor the HPA_PRIVATE flag is specified.</li> <li>• Flags: Both the HPA_PINNED and the HPA_PAGED flags are specified.</li> <li>• Flags: Neither the HPA_PINNED nor the HPA_PAGED flag is specified.</li> <li>• Keys: kernel key specified is not valid.</li> <li>• Other: Size is specified with a shared heap.</li> <li>• Other: Limit is specified with a shared heap.</li> <li>• Other: Address specified is not NULL.</li> <li>• Other: Limit &gt; size for private heap.</li> <li>• Other: Private heap size is too small (less than 8M).</li> </ul>
ENOMEM_HEAP_CREATE	Indicates insufficient memory available to complete the request.

**Related reference:**

“heap\_modify Kernel Service” on page 204  
“heap\_destroy Kernel Service”

## heap\_destroy Kernel Service

### Purpose

Removes a heap.

### Syntax

```
#include <sys/types.h>
#include <sys/malloc.h>
#include <sys/kerrno.h>
```

```
kerrno_t heap_destroy (heapattr_t heap, long flags);
```

### Parameters

Item	Description
<i>heap</i>	The heap to destroy.
<i>flags</i>	Must be zero.

### Description

This service removes a heap and its internal resources from the system. There must be no outstanding allocations when a heap is destroyed.

### Execution Environment

The *heap\_destroy* kernel service can be called from the process environment only.

### Return Values

Item	Description
EINVAL_HEAP_DESTROY	The <i>heap</i> parameter is not recognizable.
EBUSY_HEAP_DESTROY	The heap is still in use.

**Related reference:**

“heap\_create Kernel Service” on page 201  
 “heap\_modify Kernel Service”

## heap\_modify Kernel Service

### Purpose

Modifies the attributes of a heap.

### Syntax

```
#include <sys/types.h>
#include <sys/malloc.h>
#include <sys/kerrno.h>
```

```
kerrno_t heap_modify (heapattr_t heap, long command, long argument);
```

### Parameters

Item	Description
<i>heap</i>	The heap handle returned from the <b>heap_create</b> kernel service.
<i>command</i>	Specifies the operation to perform. The following values are supported:  <b>HPA_SET_LIMIT</b> Modifies the limit value of a private heap.  <b>HPA_SET_DEBUG</b> Modifies the debug level. Debug levels from 0 to 9 are supported.
<i>argument</i>	Command specific data (new limit or debug level).

### Description

The **heap\_modify** kernel service is used to alter the heap characteristics at run time.

### Execution Environment

The **heap\_modify** kernel service can be called from the process environment only with interrupts enabled.

### Return Values

Item	Description
0	Success.
EINVAL_HEAP_MODIFY	The command or the execution environment is not valid.
ERANGE_HEAP_MODIFY	Heap property is outside the supported range.

**Related reference:**

“heap\_create Kernel Service” on page 201  
 “heap\_destroy Kernel Service” on page 203

## **hkeyset\_add, hkeyset\_replace, hkeyset\_restore, or hkeyset\_get Kernel Service Purpose**

Manipulates the protection domain (page access as controlled by storage keys) in use for code execution in the kernel environment.

### **Syntax**

```
#include <sys/skeys.h>
```

```
hkeyset_t hkeyset_add ( hkeyset_t keyset );  
hkeyset_t hkeyset_replace ( hkeyset_t keyset );  
void hkeyset_restore ( hkeyset_t keyset );  
hkeyset_t hkeyset_get ( void );
```

### **Parameters**

Item	Description
<i>keyset</i>	The hardware keyset to be activated.

### **Description**

If storage protection keys are enabled, every memory page has a hardware storage protection key associated with it. A keyset is a representation of the access rights to a set of storage protection keys. To access a memory page, a hardware keyset containing the storage key associated with the memory page must be active.

The **hkeyset\_add** kernel service updates the protection domain by adding the hardware keyset specified by the *keyset* parameter to the currently addressable hardware keyset. The previous hardware keyset is returned.

The **hkeyset\_replace** kernel service updates the protection domain by loading the hardware keyset specified by the *keyset* parameter as the currently addressable storage set. The previous hardware keyset is returned.

The **hkeyset\_restore** kernel service restores a caller's hardware keyset when returning from a module entry point. It does not return any value.

The **hkeyset\_get** kernel service reads the current hardware keyset without altering it.

### **Execution Environment**

The **hkeyset\_add**, **hkeyset\_replace**, **hkeyset\_restore**, or **hkeyset\_get** kernel service can be called from either the process environment or the interrupt environment.

### **Return Values**

The **hkeyset\_add**, **hkeyset\_replace**, and **hkeyset\_get** kernel services return the *keyset* value that was active before the call. The **hkeyset\_restore** kernel service does not return any value.

#### **Related information:**

Kernel Storage Protection Keys Concepts

## **hkeyset\_restore\_userkeys Kernel Service Purpose**

Restores the previous user-memory access.

## Syntax

```
#include <sys/skeys.h>
```

```
kernno_t hkeyset_restore_userkeys (oldset)
```

```
hkeyset_t oldset;
```

## Parameters

Item	Description
<i>oldset</i>	Specifies the previous hardware keyset returned by the <b>hkeyset_update_userkeys</b> kernel service.

## Description

The **hkeyset\_restore\_userkeys** kernel service is a specialized protection gate that restores only the user-mode portion of the current hardware keyset. This is normally done by the kernel after this kernel service accesses user memory.

## Execution Environment

The **hkeyset\_restore\_userkeys** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful completion.
EINVAL_HKEYSET_RESTORE_USERKEYS	Indicates that the execution environment is not valid.

## Related reference:

“hkeyset\_update\_userkeys Kernel Service”

## hkeyset\_update\_userkeys Kernel Service

### Purpose

Establishes accessibility to user memory.

## Syntax

```
#include <sys/skeys.h>
```

```
kernno_t hkeyset_update_userkeys (oldset)
```

```
hkeyset_t *oldset;
```

## Parameters

Item	Description
<i>oldset</i>	Contains the returned previous hardware keyset. The valid parameter must be an 8-byte aligned address.

## Description

The **hkeyset\_update\_userkeys** kernel service is a specialized protection gate that alters only the user-mode portion of the current hardware keyset. The user-mode storage keys for the currently running thread is placed into the current hardware keyset. This is normally done by the kernel when this kernel service accesses user memory.

The previous hardware keyset is returned in the memory specified by the *oldset* parameter. You can use the **hkeyset\_restore\_userkeys** kernel service to remove the user accessibility when it is no longer needed.

**Important:** Kernel services such as `xmemin`, `xmemout`, `uiomove`, `copyin`, and `coypout` are the suggested ways to access user memory from the kernel. If possible, avoid using kernel code to directly access user memory.

## Execution Environment

The `hkeyset_update_userkeys` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful completion.
<code>EINVAL_HKEYSET_UPDATE_USERKEYS</code>	Indicates that the parameter or execution environment is not valid.

### Related reference:

“`hkeyset_restore_userkeys` Kernel Service” on page 205

“`xmemin` Kernel Service” on page 608

“`uiomove` Kernel Service” on page 517

## i

The following kernel services begin with the with the letter i.

### **i\_clear** Kernel Service

#### Purpose

Removes an interrupt handler.

#### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/intr.h>
```

```
void i_clear ( handler)
struct intr *handler;
```

#### Parameter

Item	Description
<i>handler</i>	Specifies the address of the interrupt handler structure passed to the <code>i_init</code> service.

#### Description

The `i_clear` service removes the interrupt handler specified by the *handler* parameter from the set of interrupt handlers that the kernel knows about. "Coding an Interrupt Handler" in *Kernel Extensions and Device Support Programming Concepts* contains a brief description of interrupt handlers.

The `i_mask` service is called by the `i_clear` service to disable the interrupt handler's bus interrupt level when this is the last interrupt handler for the bus interrupt level. The `i_clear` service removes the interrupt handler structure from the list of interrupt handlers. The kernel maintains this list for that bus interrupt level.

## Execution Environment

The `i_clear` kernel service can be called from the process environment only.

## Return Values

The `i_clear` service has no return values.

### Related reference:

“`i_init` Kernel Service” on page 217

### Related information:

I/O Kernel Services

Understanding Interrupts

## `i_disable` Kernel Service

### Purpose

Disables interrupt priorities.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/intr.h>
```

```
int i_disable ( new)
int new;
```

### Parameter

Item	Description
<i>new</i>	Specifies the new interrupt priority.

### Description

**Attention:** The `i_disable` service has two side effects that result from the replaceable and pageable nature of the kernel. First, it prevents process dispatching. Second, it ensures, within limits, that the caller's stack is in memory. Page faults that occur while the interrupt priority is not equal to `INTBASE` crash the system.

**Note:** The `i_disable` service is very similar to the standard UNIX `spl` service.

The `i_disable` service sets the interrupt priority to a more favored interrupt priority. The interrupt priority is used to control which interrupts are allowed.

A value of `INTMAX` is the most favored priority and disables all interrupts. A value of `INTBASE` is the least favored and disables only interrupts not in use. The `/usr/include/sys/intr.h` file defines valid interrupt priorities.

The interrupt priority is changed only to serialize code executing in more than one environment (that is, process and interrupt environments).

For example, a device driver typically links requests in a list while executing under the calling process. The device driver's interrupt handler typically uses this list to initiate the next request. Therefore, the device driver must serialize updating this list with device interrupts. The `i_disable` and `i_enable` services provide this ability. The `L_init` kernel service contains a brief description of interrupt handlers.

**Note:** When serializing such code in a multiprocessor-safe kernel extension, locking must be used as well as interrupt control. For this reason, new code should call the `disable_lock` kernel service instead of `i_disable`. The `disable_lock` service performs locking only on multiprocessor systems, and helps ensure that code is portable between uniprocessor and multiprocessor systems.



The **i\_disable** service must always be used with the **i\_enable** service. A routine must always return with the interrupt priority restored to the value that it had upon entry.

The **i\_mask** service can be used when a routine must disable its device across a return.

Because of these side effects, the caller of the **i\_disable** service should ensure that:

- The reference parameters are pinned.
- The code executed during the disable operation is pinned.
- The amount of stack used during the disable operation is less than 1KB.
- The called programs use less than 1KB of stack.

In general, the caller of the **i\_disable** service should also call only services that can be called by interrupt handlers. However, processes that call the **i\_disable** service can call the **e\_sleep**, **e\_wait**, **e\_sleepl**, **lockl**, and **unlockl** services as long as the event word or lockword is pinned.

The kernel's first-level interrupt handler sets the interrupt priority for an interrupt handler before calling the interrupt handler. The interrupt priority for a process is set to **INTBASE** when the process is created and is part of each process's state. The dispatcher sets the interrupt priority to the value associated with the process to be executed.

## Execution Environment

The **i\_disable** kernel service can be called from either the process or interrupt environment.

## Return Value

The **i\_disable** service returns the current interrupt priority that is subsequently used with the **i\_enable** service.

### Related reference:

“disable\_lock Kernel Service” on page 77

“i\_enable Kernel Service”

### Related information:

Understanding Interrupts

## i\_enable Kernel Service

### Purpose

Enables interrupt priorities.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/intr.h>
```

```
void i_enable ( old)
int old;
```

### Parameter

Item	Description
<i>old</i>	Specifies the interrupt priority returned by the <b>i_disable</b> service.

## Description

The **i\_enable** service restores the interrupt priority to a less-favored value. This value should be the value that was in effect before the corresponding call to the **i\_disable** service.

**Note:** When serializing a thread with an interrupt handler in a multiprocessor-safe kernel extension, locking must be used as well as interrupt control. For this reason, new code should call the **unlock\_enable** kernel service instead of **i\_enable**. The **unlock\_enable** service performs locking only on multiprocessor systems, and helps ensure that code is portable between uniprocessor and multiprocessor systems.

## Execution Environment

The **i\_enable** kernel service can be called from either the process or interrupt environment.

## Return Values

The **i\_enable** service has no return values.

### Related reference:

“i\_disable Kernel Service” on page 208

“unlock\_enable Kernel Service” on page 518

### Related information:

Understanding Interrupts

## i\_eoi Kernel Service

### Purpose

Issues an End of Interrupt (EOI) for a given handler.

### Syntax

```
int i_eoi(struct intr *handler)
```

## Description

The **i\_eoi** kernel service allows a device driver to issue an End of Interrupt (EOI) for its device explicitly. For level-triggered interrupts, after the second level interrupt handler (SLIH) has completed, the kernel issues an EOI on behalf of the device driver. For ISA (8259) edge-triggered interrupts, the kernel issues the EOI on behalf of the device driver before calling the SLIH. However, in the case of some edge-triggered interrupts (for example, PCI and PCI-E style edge-triggered interrupt), it is desirable that the device driver checks for pending work before the EOI is issued, and the driver is required to check for any additional work after the EOI is issued. The **i\_eoi** kernel service facilitates such operations and issues an EOI for an edge-triggered interrupt source. The **i\_eoi** kernel service fails if called for a level-triggered interrupt source.

## Parameters

Item	Description
<i>handler</i>	Pointer to the interrupt handler

## Execution Environment

The `i_eoi` kernel service can be called from process or interrupt environment.

## Return Values

`INTR_SUCC` if successful

`INTR_FAIL` if unsuccessful (the `INTR_EDGE` flag was not set on `i_init()`).

Virtual device drivers' interrupt services are similar to the PCI interrupt services. Interrupts are registered with a `bus_type` of `BUS_BID`. The primary difference is that the edge flag should be set for vdevices. For example:

```
Parent CuDv "bus_id" VDEVICE bus BID
Device CuAt "bus_intr_lvl" Adapter interrupt level
```

```
intr.flags |= INTR_EDGE
intr.bus_type = BUS_BID
intr.bid = Parent_CuDv.bus_id
intr.level = Device_CuAt.bus_intr_lvl
```

PCI-E interrupts are Message Signalled Interrupts, and hence, they are edge-triggered. Therefore, `INTR_EDGE` flag should be specified.

## ifa\_ifwithaddr Kernel Service

### Purpose

Locates an interface based on a complete address.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/socket.h>
#include <net/if.h>
#include <net/af.h>
```

```
struct ifaddr * ifa_ifwithaddr ( addr)
struct sockaddr *addr;
```

### Parameter

Item	Description
<i>addr</i>	Specifies a complete address.

### Description

The `ifa_ifwithaddr` kernel service is passed a complete address and locates the corresponding interface. If successful, the `ifa_ifwithaddr` service returns the `ifaddr` structure associated with that address.

## Execution Environment

The `ifa_ifwithaddr` kernel service can be called from either the process or interrupt environment.

## Return Values

If successful, the `ifa_ifwithaddr` service returns the corresponding `ifaddr` structure associated with the address it is passed. If no interface is found, the `ifa_ifwithaddr` service returns a null pointer.

## Example

To locate an interface based on a complete address, invoke the `ifa_ifwithaddr` kernel service as follows:

```
ifa_ifwithaddr((struct sockaddr *)&ipaddr);
```

### Related reference:

“`ifa_ifwithdstaddr` Kernel Service”

“`ifa_ifwithnet` Kernel Service” on page 213

### Related information:

Network Kernel Services

## `ifa_ifwithdstaddr` Kernel Service

### Purpose

Locates the point-to-point interface with a given destination address.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/socket.h>
#include <net/if.h>
```

```
struct ifaddr * ifa_ifwithdstaddr ( addr)
struct sockaddr *addr;
```

### Parameter

Item	Description
<i>addr</i>	Specifies a destination address.

### Description

The `ifa_ifwithdstaddr` kernel service searches the list of point-to-point addresses per interface and locates the connection with the destination address specified by the *addr* parameter.

### Execution Environment

The `ifa_withdstaddr` kernel service can be called from either the process or interrupt environment.

### Return Values

If successful, the `ifa_ifwithdstaddr` service returns the corresponding `ifaddr` structure associated with the point-to-point interface. If no interface is found, the `ifa_ifwithdstaddr` service returns a null pointer.

## Example

To locate the point-to-point interface with a given destination address, invoke the `ifa_ifwithdstaddr` kernel service as follows:

```
ifa_ifwithdstaddr((struct sockaddr *)&ipaddr);
```

### Related reference:

“ifa\_ifwithaddr Kernel Service” on page 211

“ifa\_ifwithnet Kernel Service”

**Related information:**

Network Kernel Services

## **ifa\_ifwithnet Kernel Service**

### **Purpose**

Locates an interface on a specific network.

### **Syntax**

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/socket.h>
#include <net/if.h>
```

```
struct ifaddr * ifa_ifwithnet ( addr)
register struct sockaddr *addr;
```

### **Parameter**

Item	Description
<i>addr</i>	Specifies the address.

### **Description**

The **ifa\_ifwithnet** kernel service locates an interface that matches the network specified by the address it is passed. If more than one interface matches, the **ifa\_ifwithnet** service returns the first interface found.

### **Execution Environment**

The **ifa\_ifwithnet** kernel service can be called from either the process or interrupt environment.

### **Return Values**

If successful, the **ifa\_ifwithnet** service returns the **ifaddr** structure of the correct interface. If no interface is found, the **ifa\_ifwithnet** service returns a null pointer.

### **Example**

To locate an interface on a specific network, invoke the **ifa\_ifwithnet** kernel service as follows:

```
ifa_ifwithnet((struct sockaddr *)&ipaddr);
```

**Related reference:**

“ifa\_ifwithaddr Kernel Service” on page 211

“ifa\_ifwithdstaddr Kernel Service” on page 212

**Related information:**

Network Kernel Services

## **if\_attach Kernel Service**

### **Purpose**

Adds a network interface to the network interface list.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <net/if.h>
```

```
if_attach ( ifp)
struct ifnet *ifp;
```

## Parameter

Item	Description
<i>ifp</i>	Points to the interface network ( <b>ifnet</b> ) structure that defines the network interface.

## Description

The **if\_attach** kernel service registers a Network Interface Driver (NID) in the network interface list.

## Execution Environment

The **if\_attach** kernel service can be called from either the process or interrupt environment.

## Return Values

The **if\_attach** kernel service has no return values.

### Related reference:

“if\_detach Kernel Service”

### Related information:

Network Kernel Services

## if\_detach Kernel Service

### Purpose

Deletes a network interface from the network interface list.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <net/if.h>
```

```
if_detach ( ifp)
struct ifnet *ifp;
```

## Parameter

Item	Description
<i>ifp</i>	Points to the interface network ( <b>ifnet</b> ) structure that describes the network interface to delete.

## Description

The **if\_detach** kernel service deletes a Network Interface Driver (NID) entry from the network interface list.

## Execution Environment

The **if\_detach** kernel service can be called from either the process or interrupt environment.

## Return Values

Item	Description
0	Indicates that the network interface was successfully deleted.
ENOENT	Indicates that the <code>if_detach</code> kernel service could not find the NID in the network interface list.

### Related reference:

“`if_attach` Kernel Service” on page 213

### Related information:

Network Kernel Services

## if\_down Kernel Service

### Purpose

Marks an interface as down.

### Syntax

```
#include <sys/types.h> #include <sys/errno.h> #include <net/if.h> void if_down ( ifp) register struct ifnet *ifp;
```

### Parameter

Item	Description
<i>ifp</i>	Specifies the <code>ifnet</code> structure associated with the interface array.

### Description

The `if_down` kernel service:

- Marks an interface as down by setting the `flags` field of the `ifnet` structure appropriately.
- Notifies the protocols of the transaction.
- Flushes the output queue.

The *ifp* parameter specifies the `ifnet` structure associated with the interface as the structure to be marked as down.

### Execution Environment

The `if_down` kernel service can be called from either the process or interrupt environment.

### Return Values

The `if_down` service has no return values.

### Example

To mark an interface as down, invoke the `if_down` kernel service as follows:

```
if_down(ifp);
```

### Related information:

Network Kernel Services

## **if\_nostat Kernel Service**

### **Purpose**

Zeroes statistical elements of the interface array in preparation for an attach operation.

### **Syntax**

```
#include <sys/types.h>
#include <sys/errno.h>
#include <net/if.h>
```

```
void if_nostat ( ifp)
struct ifnet *ifp;
```

### **Parameter**

Item	Description
<i>ifp</i>	Specifies the <b>ifnet</b> structure associated with the interface array.

### **Description**

The **if\_nostat** kernel service zeroes the statistic elements of the **ifnet** structure for the interface. The *ifp* parameter specifies the **ifnet** structure associated with the interface that is being attached. The **if\_nostat** service is called from the interface attach routine.

### **Execution Environment**

The **if\_nostat** kernel service can be called from either the process or interrupt environment.

### **Return Values**

The **if\_nostat** service has no return values.

### **Example**

To zero statistical elements of the interface array in preparation for an attach operation, invoke the **if\_nostat** kernel service as follows:

```
if_nostat(ifp);
```

#### **Related information:**

Network Kernel Services

## **ifunit Kernel Service**

### **Purpose**

Returns a pointer to the **ifnet** structure of the requested interface.

### **Syntax**

```
#include <sys/types.h>
#include <sys/errno.h>
#include <net/if.h>
```

```
struct ifnet *
ifunit ( name)
char *name;
```



## Parameter

Item	Description
<i>name</i>	Specifies the name of an interface (for example, en0).

## Description

The **ifunit** kernel service searches the list of configured interfaces for an interface specified by the *name* parameter. If a match is found, the **ifunit** service returns the address of the **ifnet** structure for that interface.

## Execution Environment

The **ifunit** kernel service can be called from either the process or interrupt environment.

## Return Values

The **ifunit** kernel service returns the address of the **ifnet** structure associated with the named interface. If the interface is not found, the service returns a null value.

## Example

To return a pointer to the **ifnet** structure of the requested interface, invoke the **ifunit** kernel service as follows:

```
ifp = ifunit("en0");
```

### Related information:

Network Kernel Services

## **i\_init** Kernel Service Purpose

Defines an interrupt handler.

## Syntax

```
#include <sys/types.h>  
#include <sys/errno.h>  
#include <sys/intr.h>
```

```
int i_init  
( handler)  
struct intr *handler;
```

## Parameter

Item	Description
<i>handler</i>	Designates the address of the pinned interrupt handler structure.

## Description

**Attention:** The interrupt handler structure must not be altered between the call to the **i\_init** service to define the interrupt handler and the call to the **i\_clear** service to remove the interrupt handler. The structure must also stay pinned. If this structure is altered at those times, a kernel panic may result.

The **i\_init** service allows device drivers to define an interrupt handler to the kernel. The interrupt handler **intr** structure pointed to by the *handler* parameter describes the interrupt handler. The caller of the **i\_init** service must initialize all the fields in the **intr** structure. The `/usr/include/sys/intr.h` file defines these fields and their valid values.

The **i\_init** service enables interrupts by linking the interrupt handler structure to the end of the list of interrupt handlers defined for that bus level. If this is the first interrupt handler for the specified bus interrupt level, the **i\_init** service enables the bus interrupt level by calling the **i\_unmask** service.

The interrupt handler can be called before the **i\_init** service returns if the following two conditions are met:

- The caller of the **i\_init** service is executing at a lower interrupt priority than the one defined for the interrupt.
- An interrupt for the device or another device on the same bus interrupt level is already pending.

On multiprocessor systems, all interrupt handlers defined with the **i\_init** kernel service run by default on the first processor started when the system was booted. This ensures compatibility with uniprocessor interrupt handlers. If the interrupt handler being defined has been designed to be multiprocessor-safe, or is an EPOW (Early Power-Off Warning) or off-level interrupt handler, set the **INTR\_MPSAFE** flag in the `flags` field of the **intr** structure passed to the **i\_init** kernel service. The interrupt handler will then run on any available processor.

### Coding an Interrupt Handler

The kernel calls the interrupt handler when an enabled interrupt occurs on that bus interrupt level. The interrupt handler is responsible for determining if the interrupt is from its own device and processing the interrupt. The interface to the interrupt handler is as follows:

```
int interrupt_handler (handler) struct intr *handler;
```

The *handler* parameter points to the same interrupt handler structure specified in the call to the **i\_init** kernel service. The device driver can pass additional parameters to its interrupt handler by declaring the interrupt handler structure to be part of a larger structure that contains these parameters.

The interrupt handler can return one of two return values. A value of **INTR\_SUCC** indicates that the interrupt handler processed the interrupt and reset the interrupting device. A value of **INTR\_FAIL** indicates that the interrupt was not from this interrupt handler's device.

### Registering Early Power-Off Warning (EPOW) Routines

The **i\_init** kernel service can also be used to register an EPOW (Early Power-Off Warning) notification routine.

The return value from the EPOW interrupt handler should be **INTR\_SUCC**, which indicates that the interrupt was successfully handled. All registered EPOW interrupt handlers are called when an EPOW interrupt is indicated.

### Execution Environment

The **i\_init** kernel service can be called from the process environment only.

### Return Values

Item	Description
INTR_SUCC	Indicates a successful completion.
INTR_FAIL	Indicates an unsuccessful completion. The <code>i_init</code> service did not define the interrupt handler.

An unsuccessful completion occurs when there is a conflict between a shared and a nonshared bus interrupt level. An unsuccessful completion also occurs when more than one interrupt priority is assigned to a bus interrupt level.

#### Related information:

Understanding Interrupts

I/O Kernel Services

## **i\_mask Kernel Service**

### **Purpose**

Disables a bus interrupt level.

### **Syntax**

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/intr.h>
```

```
void i_mask ( handler)
struct intr *handler;
```

### **Parameter**

Item	Description
<i>handler</i>	Specifies the address of the interrupt handler structure that was passed to the <code>i_init</code> service.

### **Description**

The `i_mask` service disables the bus interrupt level specified by the *handler* parameter.

The `i_disable` and `i_enable` services are used to serialize the execution of various device driver routines with their device interrupts.

The `i_init` and `i_clear` services use the `i_mask` and `i_unmask` services internally to configure bus interrupt levels.

Device drivers can use the `i_disable`, `i_enable`, `i_mask`, and `i_unmask` services when they must perform off-level processing with their device interrupts disabled. Device drivers also use these services to allow process execution when their device interrupts are disabled.

### **Execution Environment**

The `i_mask` kernel service can be called from either the process or interrupt environment.

### **Return Values**

The `i_mask` service has no return values.

#### **Related reference:**

“`i_unmask Kernel Service`” on page 239

#### **Related information:**

Understanding Interrupts

I/O Kernel Services

## | **in\_localaddr Kernel Service**

### | **Purpose**

| Determine whether an IPv4 address is on the local network

### | **Syntax**

| #include <arpa/inet.h>

| **int in\_localaddr** (*struct in\_addr in*)

### | **Parameters**

| **in** Specifies the IPv4 address

### | **Description**

| Indicates that the IPv4 address in is for a local host (one to which we have a connection). If subnetsarelocal is true, this includes other subnets of the local net. Otherwise, it includes only the directly-connected (sub)nets.

### | **Execution Environment**

| The **in\_localaddr** kernel service can be called from the process environment only.

### | **Return Values**

| **1** Indicates that the internet address is for a local host (one to which we have a connection). If subnetsarelocal is true, this includes other subnets of the local net. Otherwise, it includes only the directly-connected (sub)nets.

## | **init\_heap Kernel Service**

### **Purpose**

Initializes a new heap to be used with kernel memory management services.

### **Syntax**

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/xmalloc.h>
#include <sys/malloc.h>
```

```
heapaddr_t init_heap ( area, size, heapp)
caddr_t area;
int size;
heapaddr_t *heapp;
```

### **Parameters**

Item	Description
<i>area</i>	Specifies the virtual memory address used to define the starting memory area for the heap. This address must be page-aligned.
<i>size</i>	Specifies the size of the heap in bytes. This value must be an integral number of system pages.
<i>heapp</i>	Points to the external heap descriptor. This must have a null value. The base kernel uses this field is used to specify special heap characteristics that are unavailable to kernel extensions.

## Description

The `init_heap` kernel service is most commonly used by a kernel process to initialize and manage an area of virtual memory as a private heap. Once this service creates a private heap, the returned `heapaddr_t` value can be used with the `xmalloc` or `xmfree` service to allocate or deallocate memory from the private heap. Heaps can be created within other heaps, a kernel process private region, or even on a stack.

Few kernel extensions ever require the `init_heap` service because the exported global `kernel_heap` and `pinned_heap` are normally used for memory allocation within the kernel. However, kernel processes can use the `init_heap` service to create private nonglobal heaps within their process private region for controlling kernel access to the heap and possibly for performance considerations.

## Execution Environment

The `init_heap` kernel service can be called from the process environment only.

### Related reference:

“`xmalloc` Kernel Service” on page 598

### Related information:

Memory Kernel Services

Using Kernel Processes

## initp Kernel Service

### Purpose

Changes the state of a kernel process from idle to ready.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
int initp
(pid, func, init_parms,
parms_length, name)
pid_t pid;
void (func) (int
flag, void* init_parms, int parms_length );
void * init_parms;
int parms_length;
char * name;
```

### Parameters

Item	Description
<i>pid</i>	Specifies the process identifier of the process to be initialized.
<i>func</i>	Specifies the process's initialization routine.
<i>init_parm</i>	Specifies the pointer to the initialization parameters.
<i>parms_length</i>	Specifies the length of the initialization parameters.
<i>name</i>	Specifies the process name.

## Description

The **initp** kernel service completes the transition of a kernel process from idle to ready. The idle state for a process is represented by **p\_status == SIDL**. Before calling the **initp** service, the **creatp** service is called to create the process. The **creatp** service allocates and initializes a process table entry.

The **initp** service creates and initializes the process-private segment. The process is marked as a kernel process by a bit set in the **p\_flag** field in the process table entry. This bit, the SKPROC bit, signifies that the process is a kernel process.

The process calling the **initp** service to initialize a newly created process must be the same process that called the **creatp** service to create the new process.

"Using Kernel Processes" in *Kernel Extensions and Device Support Programming Concepts* further explains how the **initp** kernel service completes the initialization process begun by the **creatp** service.

The *pid* parameter identifies the process to be initialized. It must be valid and identify a process in the SIDL (idle) state.

The *name* parameter points to a character string that names the process. The leading characters of this string are copied to the user structure. The number of characters copied is implementation-dependent, but at least four are always copied.

The *func* parameter indicates the main entry point of the process. The new process is made ready to run this function. If the *init\_parms* parameter is not null, it points to data passed to this routine. The parameter structure must be agreed upon between the initializing and initialized process. The **initp** service copies the data specified by the *init\_parm* parameter (with the exact number of bytes specified by the *parms\_length* parameter) of data to the new process's stack.

## Execution Environment

The **initp** kernel service can be called from the process environment only.

## Example

To initialize the kernel process running the function *main\_kproc*, enter:

```
{
.
.
.
pid = creatp();
initp(pid, main_kproc, &node_num, sizeof(int), "tkproc");
.
.
}
void
main_kproc(int flag, void* init_parms, int parms_length)
{
.
.
}
```

```

    .
    int i;
    i = *( (int *)init_parms );
    .
    .
}

```

## Return Values

Item	Description
0	Indicates a successful operation.
ENODEV	The process could not be scheduled because it has a processor attachment that does not contain any available processors. This can be caused by Dynamic Processor Deallocation.
ENOMEM	Indicates that there was insufficient memory to initialize the process.
EINVAL	Indicates an <i>pid</i> parameter that was not valid.

### Related reference:

“creatp Kernel Service” on page 54

### Related information:

Process and Exception Management Kernel Services

Dynamic logical partitioning

## initp Kernel Service func Subroutine

### Purpose

Directs the process initialization routine.

### Syntax

```

#include <sys/types.h>
#include <sys/errno.h>

```

```

void func (flag, init_parms, parms_length)
int flag;
void * init_parms;
int parms_length;

```

### Parameters

Item	Description
<i>func</i>	Specifies the process's initialization routine.
<i>flag</i>	Has a 0 value if this subroutine is executed as a result of initializing a process with the <b>initp</b> service.
<i>init_parms</i>	Specifies the pointer to the initialization parameters.
<i>parms_length</i>	Specifies the length of the initialization parameters.

### Related reference:

“initp Kernel Service” on page 221

### Related information:

Process and Exception Management Kernel Services

## io\_map Kernel Service

### Purpose

Attach to an I/O mapping

## Syntax

```
#include <sys/adspace.h>
```

```
void * io_map (io_handle)  
io_handle_t io_handle;
```

## Description

The **io\_map** kernel service sets up addressability to the I/O address space defined by the **io\_handle\_t** structure. It returns an effective address representing the start of the mapped region.

The **io\_map** kernel service is a replacement call for the **iomem\_att** kernel service, which is deprecated on AIX 6.1. However, the **io\_map** kernel service might replace multiple **iomem\_att** calls depending on the device, the driver, and whether multiple regions were mapped into a single virtual segment. Like the **iomem\_att** kernel service, this service does not return any kind of failure. If something goes wrong, the system crashes.

There is a major difference between **io\_map** and **iomem\_att**. **iomem\_att** took an **io\_map** structure containing a bus address and returned a fully qualified effective address with any byte offset from the bus address preserved and computed into the returned effective address. The **io\_map** kernel service always returns a segment-aligned effective address representing the beginning of the I/O segment corresponding to **io\_handle\_t**. Manipulation of page and byte offsets within the segment are responsibilities of the device driver.

The **io\_map** kernel service is subject to nesting rules regarding the number of attaches allowed. A total system number of active temporary attaches is 4. However, it is recommended that no more than one active attach be owned by a driver calling the interrupt or DMA kernel services. It is also recommended that no active attaches be owned by a driver when calling other kernel services.

## Parameters

Item	Description
<i>io_handle</i>	Received on a prior successful call to <code>io_map_init</code> . Describes the I/O space to attach to.

## Execution Environment

The **io\_map** kernel service can be called from the process or interrupt environment.

## Return Values

The **io\_map** kernel service returns a segment-aligned effective address to access the I/O address spaces.

### Related reference:

“`io_map_init` Kernel Service” on page 225

“`io_unmap` Kernel Service” on page 226

### Related information:

Programmed I/O (PIO) Kernel Services

## **io\_map\_clear** Kernel Service

### Purpose

Removes an I/O mapping segment.



## Syntax

```
#include <sys/adspc.h>
```

```
void io_map_clear (io_handle)
io_handle_t io_handle;
```

## Description

This service destroys all mappings defined by the *io\_handle\_t* parameter.

There should be no active mappings (outstanding **io\_map** calls) to this handle when **io\_map\_clear** is called. The segment previously created by an **io\_map\_init** call or multiple **io\_map\_init** calls, is deleted.

## Parameters

Item	Description
<i>io_handle</i>	Received on a prior successful call to <b>io_map_init</b> . Describes the I/O space to be removed.

## Execution Environment

The **io\_map\_clear** kernel service can be called from the process environment only.

### Related reference:

“io\_map Kernel Service” on page 223

“io\_unmap Kernel Service” on page 226

### Related information:

Programmed I/O (PIO) Kernel Services

## io\_map\_init Kernel Service

### Purpose

Creates and initializes an I/O mapping segment.

## Syntax

```
#include <sys/adspc.h>
```

```
#include <sys/vm_types.h>
```

```
io_handle_t io_map_init (io_map_ptr, page_offset, io_handle)
```

```
struct io_map *io_map_ptr;
```

```
vpn_t page_offset;
```

```
io_handle_t io_handle;
```

```
struct io_map {
    int key;                /* structure version number */
    int flags;              /* flags for mapping */
    int32long64_t size;     /* size of address space needed */
    int bid;                /* bus ID */
    long long busaddr;     /* bus address */
};
```

## Description

The **io\_map\_init** kernel service will create a segment to establish a cache-inhibited virtual-to-real translation for the bus address region defined by the contents of the **io\_map** struct. The *flags* parameter of the **io\_map** structure can be used to customize the mapping such as making the region read-only, using the **IOM\_RDONLY** flag.

The `io_map_init` kernel service returns a handle of an opaque type `io_handle_t` to be used on future `io_map` or `io_unmap` calls. All services that use the `io_handle` returned by `io_map_init` must use the handle from the most recent call. Using an old handle is a programming error.

The `vpn_t` type parameter represents the virtual page number offset to allow the caller to specify where, in the virtual segment, to map this region. The offset must not conflict with a previous mapping in the segment. The caller should map the most frequently accessed and performance critical I/O region at `vpn_t` offset 0 into the segment. This is due to the fact that the subsequent `io_map` calls using this `io_handle` will return an effective address representing the start of the segment (that is, page offset 0). The device driver is responsible for managing various offsets into the segment. A single bus memory address page can be mapped multiple times at different `vpn_t` offsets within the segment.

The `io_handle_t` parameter is useful when the caller wants to append a new mapping to an existing segment. For the initial creation of a new I/O segment, this parameter must be NULL. For appended mappings to the same segment, this parameter is the `io_handle_t` returned from the last successful `io_map_init` call. If the mapping fails for any reason (offset conflicts with prior mapping, or no more room in the segment), NULL is returned. In this case, the previous `io_handle_t` is still valid. If successful, the `io_handle_t` returned should be used on all future calls. In this way, a device driver can manage multiple I/O address spaces of a single adapter within a single virtual address segment, requiring the driver to do only a single attach, `io_map`, to gain addressability to all of the mappings.

## Parameters

Item	Description
<code>io_map_ptr</code>	Pointer to <code>io_map</code> structure describing the address region to map.
<code>page_offset</code>	Page offset at which to map the specified region into the virtual address segment.
<code>io_handle</code>	For the first call, this parameter should be NULL. When adding to an existing mapping, this parameter is the <code>io_handle</code> received on a prior successful call to <code>io_map_init</code> .

## Execution Environment

The `io_map_init` kernel service can be called from the process environment only.

## Return Values

Item	Description
<code>io_handle_t</code>	An opaque handle to the mapped I/O segment in the virtual memory that must be used in subsequent calls to this service.
NULL	Failed to create or append mapping.

### Related reference:

“`io_map_clear` Kernel Service” on page 224

“`io_unmap` Kernel Service”

### Related information:

Programmed I/O (PIO) Kernel Services

## `io_unmap` Kernel Service Purpose

Detach from an I/O mapping

## Syntax

```
#include <sys/adspace.h>
```

```
void io_unmap (eaddr)  
void *eaddr;
```

## Description

The **io\_unmap** kernel service removes addressability to the I/O address space defined by the *eaddr* parameter. There must be a valid active mapping from a previous **io\_map** call for this effective address. The *eaddr* parameter can be any valid effective address within the segment, and it does not have to be exactly the same as the address returned by **io\_map**.

The **io\_unmap** kernel service is a replacement call for the **iomem\_det** kernel service, which is deprecated on AIX 6.1. However, the **io\_unmap** kernel service might replace multiple **iomem\_det** calls depending on the device, the driver, and whether multiple regions were mapped into a single virtual segment using the **io\_map\_init** kernel service.

## Parameters

Item	Description
<i>eaddr</i>	Received on a prior successful call to <b>io_map</b> . Effective address for the I/O space to detach from.

## Execution Environment

The **io\_unmap** kernel service can be called from the process or interrupt environment.

### Related reference:

“**io\_map\_clear** Kernel Service” on page 224

“**io\_map** Kernel Service” on page 223

### Related information:

Programmed I/O (PIO) Kernel Services

## iodone Kernel Service

### Purpose

Performs block I/O completion processing.

## Syntax

```
#include <sys/types.h>  
#include <sys/errno.h>  
#include <sys/buf.h>
```

```
void iodone ( bp)  
struct buf *bp;
```

## Parameter

Item	Description
<i>bp</i>	Specifies the address of the <b>buf</b> structure for the buffer whose I/O has completed.  On a platform that supports storage keys, the passed in <i>bp</i> parameter must be in the <b>KKEY_PUBLIC</b> or <b>KKEY_BLOCK_DEV</b> protection domain.

## Description

A device driver calls the **iodone** kernel service when a block I/O request is complete. The device driver must not reference or alter the buffer header or buffer after calling the **iodone** service.

The **iodone** service takes one of two actions, depending on the current interrupt level. Either it invokes the caller's individual **iodone** routine directly, or it schedules I/O completion processing for the buffer to be performed off-level, at the **INTIODONE** interrupt level. The interrupt handler for this level then calls the **iodone** routine for the individual device driver. In either case, the individual **iodone** routine is defined by the **b\_iodone** buffer header field in the buffer header. This **iodone** routine is set up by the caller of the device's strategy routine.

For example, the file I/O system calls set up a routine that performs buffered I/O completion processing. The **uphysio** service sets up a routine that performs raw I/O completion processing. Similarly, the pager sets up a routine that performs page-fault completion processing.

## Setting up an iodone Routine

Under certain circumstances, a device driver can set up an **iodone** routine. For example, the logical volume device driver can follow this procedure:

1. Take a request for a logical volume.
2. Allocate a buffer header.
3. Convert the logical volume request into a physical volume request.
4. Update the allocated buffer header with the information about the physical volume request. This includes setting the **b\_iodone** buffer header field to the address of the individual **iodone** routine.
5. Call the physical volume device driver strategy routine.

Here, the caller of the logical volume strategy routine has set up an **iodone** routine that is started when the logical volume request is complete. The logical volume strategy routine in turn sets up an **iodone** routine that is invoked when the physical volume request is complete.

The key point of this example is that only the caller of a strategy routine can set up an **iodone** routine and even then, this can only be done while setting up the request in the buffer header.

The interface for the **iodone** routine is identical to the interface to the **iodone** service.

## Execution Environment

The **iodone** kernel service can be called from either the process or interrupt environment.

## Return Values

The **iodone** service has no return values.

### Related reference:

"iowait Kernel Service" on page 233

"buf Structure" on page 615

### Related information:

I/O Kernel Services

## **iostadd Kernel Service**

### **Purpose**

Registers an I/O statistics structure that is used for updating I/O statistics reported by the **iostat** subroutine.

### **Syntax**

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/iostat.h>
#include <sys/devinfo.h>

int iostadd ( devtype, devstatp)
int devtype;
union {
    struct ttystat *ttystp;
    struct dkstat *dkstp;
} devstatp;
```

### **Description**

The **iostadd** kernel service is used to register the I/O statistics structure that is required to maintain statistics on a device. The **iostadd** service is typically called by a tty, disk, or CD-ROM device driver to provide the statistical information that is used by the **iostat** subroutine. The **iostat** subroutine displays statistic information for tty and disk devices on the system. The **iostadd** service must be used once for each configured device.

The **iostadd** kernel service and the **dkstat** structure support Multi-Path I/O (MPIO). For an MPIO device, the anchor is the disk's **dkstat** structure. This anchor must be the first **dkstat** structure that is registered by using the **iostadd** kernel service. Any path **dkstat** structures that are registered later must reference the address of the anchor **dkstat** (disk) structure in the `dkstat.dk_mpio_anchor` field.

For tty devices, the *devtype* parameter has a value of **DD\_tty**. In this case, the **iostadd** service uses the *devstatp* parameter to return a pointer to a **ttystat** structure.

For disk or CD-ROM devices with a *devtype* value of **DD\_DISK** or **DD\_CD-ROM**, the caller must provide a pinned and initialized **dkstat** structure as an input parameter. This structure is pointed to by the *devstatp* parameter on entry to the **iostadd** kernel service.

If the device driver support for a device is terminated, the **dkstat** or **ttystat** structure that is registered with the **iostadd** kernel service must be deregistered by calling the **iostdel** kernel service.

### **I/O Statistics Structures**

The **iostadd** kernel service uses two structures that are found in the **usr/include/sys/iostat.h** file: the **ttystat** structure and the **dkstat** structure.

The **ttystat** structure contains the following fields:

Field	Description
rawinch	Count of raw characters that are received by the tty device
caninch	Count of canonical characters that are generated from canonical processing
outch	Count of the characters output to a tty device

The second structure that is used by the **iostadd** kernel service is the **dkstat** structure, which contains information about disk devices. This structure contains the following fields:

Field	Description
diskname	32-character string name for the disk's logical device
dknextp	Pointer to the next <b>dkstat</b> structure in the chain
dk_status	Disk entry-status flags
dk_time	Time the disk is active
dk_bsize	Number of bytes in a block
dk_xfers	Number of transfers to or from the disk
dk_rblks	Number of blocks that are read from the disk
dk_wblks	Number of blocks that are written to the disk
dk_seeks	Number of seek operations for disks
dk_version	Version of the <b>dkstat</b> structure
dk_q_depth	Queue depth
dk_mpio_anchor	Pointer to the path data anchors (disk)
dk_mpio_next_path	Pointer to the next path <b>dkstat</b> structure in the chain
dk_mpio_path_id	Path ID

## tty Device Driver Support

The `rawinch` field in the **ttystat** structure must be incremented by the number of characters that are received by the tty device. The `caninch` field in the **ttystat** structure must be incremented by the number of input characters that are generated from canonical processing. The `outch` field is increased by the number of characters output to tty devices. These fields must be incremented by the device driver, but never be cleared.

## Disk Device Driver Support

A disk device driver must perform these four tasks:

- Allocate and pin a **dkstat** structure during device initialization.
- Update the `dkstat.diskname` field with the device's logical name.
- Update the `dkstat.dk_bsize` field with the number of bytes in a block on the device.
- Set all other fields in the structure to 0.

If a disk device driver supports MPIO, it must perform the following tasks:

- Allocate and pin a **dkstat** structure during device initialization.
- Update the `dkstat.diskname` field with the device's logical name.
- Update the `dkstat.dk_bsize` field with the number of bytes in a block on the device.
- Set the value of `dkstat.dk_version` to **dk\_qd\_mpio\_magic**.
- Set the value of `dkstat.dk_mpio_anchor` to 0 if the **dkstat** added structure is the disk.
- Set the value of `dkstat.dk_mpio_anchor` to the address of the path's anchor (disk) **dkstat** structure, and set `dkstat.dk_mpio_path_id` to the path's ID if the **dkstat** added structure is a path.
- Set all other fields to 0.

If the device supports discrete seek commands, the `dkstat.dk_xrate` field in the structure must be set to the transfer rate capability of the device (KB/sec). The device's **dkstat** structure must then be registered by using the **iostadd** kernel service.

During drive operation update, the `dkstat.dk_status` field must show the busy or non-busy state of the device. It can be done by setting and resetting the `IOST_DK_BUSY` flag. The `dkstat.dk_xfers` field must be incremented for each transfer initiated to or from the device. The `dkstat.dk_rblks` and `dkstat.dk_wblks` fields must be incremented by the number of blocks that are read or written.

If the device supports discrete seek commands, the `dkstat.dk_seek` field must be incremented by the number of seek commands that are sent to the device. If the device does not support discrete seek commands, both the `dkstat.dk_seek` and `dkstat.dk_xrate` fields must be left with a value of 0.

The base kernel updates the `dkstat.dk_nextp` and `dkstat.dk_time` fields. They must not be modified by the device driver after initialization. For MPIO devices, the base kernel also updates the `dkstat.dk_mpio_next_path` field.

**Note:** The same `dkstat` structure must not be registered more than once.

In addition to basic tasks, a disk driver must perform the following tasks before it calls the `iostadd` kernel service if the driver supports the `-D` option of the `iostat` command:

- Set the value of `dkstat.dk_version` to `dk_qd_service2_magic`.
- Set the `dkstat.ident.adapter` field to the adapter name if the driver does not support MPIO.

During I/O operations, the driver must perform the following tasks:

- Increase the `dkstat.__dk_rxfers` field for each read transfer.
- Update the `dkstat.dk_q_depth` field with the number of I/O requests which are in progress.
- Increase the `dkstat.dk_q_full` field when the number of I/O requests which are in progress reaches the maximum queue depth.
- Increase the `dkstat.dk_rserv` field by the service time which is the delta-time base value between when the `devstrat` kernel service sends a read request to the adapter driver and when the `iodone` kernel service returns the request from the adapter driver.
- Increase the `dkstat.dk_rtimeout` field when the driver tries a failed read request again.
- Increase the `dkstat.dk_rfailed` field when the driver returns a failed read request as an error.
- Update the `dkstat.dk_min_rserv` field with the minimum service time for a read request.
- Update the `dkstat.dk_max_rserv` field with the maximum service time for a read request.
- Increase the `dkstat.dk_wserv` field by the service time which is the delta-time base value between when the `devstrat` kernel service sends a write request to the adapter driver and when the `iodone` kernel service returns the request from the adapter driver.
- Increase the `dkstat.dk_wtimeout` field when the driver tries a failed write request again.
- Increase the `dkstat.dk_wfailed` field when the driver returns a failed write request as an error.
- Update the `dkstat.dk_min_wserv` field with the minimum service time for a write request.
- Update the `dkstat.dk_max_wserv` field with the maximum service time for a write request.
- Increase and decrease the `dkstat.dk_wq_depth` field when the driver enqueues and dequeues an I/O request.
- Increase the `dkstat.dk_wq_time` field by the wait time which is the delta-time base value between when the driver enqueues an I/O request and when the `devstrat` kernel service sends the request to the adapter driver.
- Update the `dkstat.dk_wq_min_time` field with the minimum wait time.
- Update the `dkstat.dk_wq_max_time` field with the maximum wait time.

## Parameters

Item	Description
<i>devtype</i>	Specifies the type of device for which I/O statistics are kept. The various device types are defined in the <code>/usr/include/sys/devinfo.h</code> file. Currently, I/O statistics are only kept for disks, CD-ROMs, and tty devices. Possible values for this parameter are:  <b>DD_DISK</b> For disks  <b>DD_CD-ROM</b> For CD-ROMs  <b>DD_TTY</b> For tty devices
<i>devstatp</i>	Points to an I/O statistics structure for the device type that is specified by the <i>devtype</i> parameter. For a <i>devtype</i> parameter of <b>DD_tty</b> , the address of a pinned <b>ttystat</b> structure is returned. For a <i>devtype</i> parameter of <b>DD_DISK</b> or <b>DD_CD-ROM</b> , the parameter is an input parameter that points to a <b>dkstat</b> structure previously allocated by the caller.  On a platform that supports storage keys, the passed in <i>devstatp</i> parameter must be in the <b>KKEY_PUBLIC</b> or <b>KKEY_BLOCK_DEV</b> protection domain.

## Execution Environment

The **iostadd** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates that no error is detected.
EINVAL	Indicates that the <i>devtype</i> parameter specified a device type that is not valid. For MPIO devices, indicates that an anchor for a path <b>dkstat</b> structure was not found.

### Related reference:

“iostdel Kernel Service”

### Related information:

iostat subroutine

Kernel Extension and Device Driver Management Kernel Services

## iostdel Kernel Service

### Purpose

Removes the registration of an I/O statistics structure that is used for maintaining I/O statistics on a particular device.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/iostat.h>
```

```
void iostdel ( devstatp)
union {
    struct ttystat *ttystp;
    struct dkstat *dkstp;
} devstatp;
```



## Description

The **iostdel** kernel service removes the registration of an I/O statistics structure for a terminating device. The device's **ttystat** or **dkstat** structure must be previously registered by using the **iostadd** kernel service. Following a return from the **iostdel** service, the **iostat** command no longer displays statistics for the device that is terminated.

The **iostdel** kernel service supports Multi-Path I/O (MPIO). For an MPIO device, the anchor is the disk's **dkstat** structure. An anchor (disk) might have several paths that are associated with it. Each of these paths can have a **dkstat** structure that is registered by using the **iostadd** kernel service. The semantics for unregistering a **dkstat** structure for an MPIO device are more restrictive than for a non-MPIO device. All paths must unregister before the anchor (disk) is unregistered. If the anchor (disk) **dkstat** structure is unregistered before all of the paths that are associated with it are unregistered, the **iostdel** kernel service removes the registration of the anchor (disk) **dkstat** structure and all remaining registered paths.

## Parameters

Item	Description
<i>devstatp</i>	Points to an I/O statistics structure previously registered by using the <b>iostadd</b> kernel service.  On a platform that supports storage keys, the passed in <i>devstatp</i> parameter must be in the <b>KKEY_PUBLIC</b> or <b>KKEY_BLOCK_DEV</b> protection domain.

## Execution Environment

The **iostdel** kernel service can be called from the process environment only.

## Return Values

The **iostdel** service has no return values.

### Related reference:

“iostadd Kernel Service” on page 229

### Related information:

iostat command

Kernel Extension and Device Driver Management Kernel Services

## iowait Kernel Service

### Purpose

Waits for block I/O completion.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/buf.h>
```

```
int iowait ( bp)
struct buf *bp;
```

### Parameter

Item	Description
------	-------------

<i>bp</i>	Specifies the address of the <b>buf</b> structure for the buffer with in-process I/O.
-----------	---

On a platform that supports storage keys, the passed in *bp* parameter must be in the **KKEY\_PUBLIC** or **KKEY\_BLOCK\_DEV** protection domain.

## Description

The **iodone** kernel service causes a process to wait until the I/O is complete for the buffer specified by the *bp* parameter. Only the caller of the strategy routine can call the **iodone** service. The **B\_ASYNC** bit in the buffer's *b\_flags* field should not be set.

The **iodone** kernel service must be called when the block I/O transfer is complete. The **buf** structure pointed to by the *bp* parameter must specify an **iodone** routine. This routine is called by the **iodone** interrupt handler in response to the call to the **iodone** kernel service. This **iodone** routine must call the **e\_wakeup** service with the *bp->b\_events* field as the event. This action awakens all processes waiting on I/O completion for the **buf** structure using the **iodone** service.

## Execution Environment

The **iodone** kernel service can be called from the process environment only.

## Return Values

The **iodone** service uses the **geterror** service to determine which of the following values to return:

Item	Description
0	Indicates that I/O was successful on this buffer.
EIO	Indicates that an I/O error has occurred.
<i>b_error</i> value	Indicates that an I/O error has occurred on the buffer.

### Related reference:

“geterror Kernel Service” on page 187

“iodone Kernel Service” on page 227

“buf Structure” on page 615

## ip\_fltr\_in\_hook, ip\_fltr\_out\_hook, ipsec\_decap\_hook, inbound\_fw, outbound\_fw Kernel Service

### Purpose

Contains hooks for IP filtering.

### Syntax

```
#define FIREWALL_OK      0 /* Accept IP packet          */
#define FIREWALL_NOTOK  1 /* Drop IP packet           */
#define FIREWALL_OK_NOTSEC 2 /* Accept non-encapsulated IP packet
                        (ipsec_decap_hook only) */

#include <sys/mbuf.h>
#include <net/if.h>

int (*ip_fltr_in_hook)(struct mbuf **pkt, void **arg)

int (*ipsec_decap_hook)(struct mbuf **pkt, void **arg)

int (*ip_fltr_out_hook)(struct ifnet *ifp, struct mbuf **pkt,
int flags)
```

```

#include <sys/types.h>

#include <sys/mbuf.h>

#include <netinet/ip_var.h>

void (*inbound_fw)(struct ifnet *ifp, struct mbuf *m,
inbound_fw_args_t *args)

void ipintr_noqueue_post_fw(struct ifnet *ifp, struct mbuf *m,
inbound_fw_args_t *args)

inbound_fw_args_t *inbound_fw_save_args(inbound_fw_args_t *args)

int (*outbound_fw)(struct ifnet *ifp, struct mbuf *
m0, outbound_fw_args_t *args)

int ip_output_post_fw( struct ifnet *ifp, struct mbuf *m0,
outbound_fw_args_t *args)

outbound_fw_args_t *outbound_fw_save_args(outbound_fw_args_t *args)

```

## Parameters

Item	Description
<i>pkt</i>	Points to the mbuf chain containing the IP packet to be received ( <b>ip_fltr_in_hook</b> , <b>ipsec_decap_hook</b> ) or transmitted ( <b>ip_fltr_out_hook</b> ). The <i>pkt</i> parameter may be examined and/or changed in any of the three hook functions.
<i>arg</i>	Is the address of a pointer to <i>void</i> that is locally defined in the function where <b>ip_fltr_in_hook</b> and <b>ipsec_decap_hook</b> are called. The <i>arg</i> parameter is initially set to NULL, but the address of this pointer is passed to the two hook functions, <b>ip_fltr_in_hook</b> and <b>ipsec_decap_hook</b> . The <i>arg</i> parameter may be set by either of these functions, thereby allowing a void pointer to be shared between them.
<i>ifp</i>	Is the outgoing interface on which the IP packet will be transmitted for the <b>ip_fltr_out_hook</b> function.
<i>flags</i>	Indicates the ip_output flags passed by a transport layer protocol. Valid flags are currently defined in the <code>/usr/include/netinet/ip_var.h</code> files. See the Flags section below.

## Description

These routines provide kernel-level hooks for IP packet filtering enabling IP packets to be selectively accepted, rejected, or modified during reception, transmission, and decapsulation. These hooks are initially NULL, but are exported by the netinet kernel extension and will be invoked if assigned non-NULL values.

The **ip\_fltr\_in\_hook** routine is used to filter incoming IP packets, the **ip\_fltr\_out\_hook** routine filters outgoing IP packets, and the **ipsec\_decap\_hook** routine filters incoming encapsulated IP packets.

The **ip\_fltr\_in\_hook** function is invoked for every IP packet received by the host, whether addressed directly to this host or not. It is called after verifying the integrity and consistency of the IP packet. The function is free to examine or change the IP packet (*pkt*) or the pointer shared with **ipsec\_decap\_hook** (*arg*). The return value of the **ip\_fltr\_in\_hook** indicates whether *pkt* should be accepted or dropped. The return values are described in Expected Return Values below. If *pkt* is accepted (a return value of **FIREWALL\_OK**) and it is addressed directly to the host, the **ipsec\_decap\_hook** function is invoked next. If *pkt* is accepted, but is not directly addressed to the host, it is forwarded if IP forwarding is enabled. If **ip\_fltr\_in\_hook** indicates *pkt* should be dropped (a return value of **FIREWALL\_NOTOK**), it is neither delivered nor forwarded.

The **ipsec\_decap\_hook** function is called after reassembly of any IP fragments (the **ip\_fltr\_in\_hook** function will have examined each of the IP fragments) and is invoked only for IP packets that are directly addressed to the host. The **ipsec\_decap\_hook** function is free to examine or change the IP packet (*pkt*) or the pointer shared with **ipsec\_decap\_hook** (*arg*). The hook function should perform decapsulation if necessary, back into *pkt* and return the proper status so that the IP packet can be processed appropriately.

See the Expected Return Values section below. For acceptable encapsulated IP packets (a return value of **FIREWALL\_OK**), the decapsulated packet is processed again by jumping to the beginning of the IP input processing loop. Consequently, the decapsulated IP packet will be examined first by **ip\_fltr\_in\_hook** and, if addressed to the host, by **ipsec\_decap\_hook**. For acceptable non-encapsulated IP packets (a return value of **FIREWALL\_OK\_NOTSEC**), IP packet delivery simply continues and *pkt* is processed by the transport layer. A return value of **FIREWALL\_NOTOK** indicates that *pkt* should be dropped.

The **ip\_fltr\_out\_hook** function is called for every IP packet to be transmitted, provided the outgoing IP packet's destination IP address is NOT an IP multicast address. If it is, it is sent immediately, bypassing the **ip\_fltr\_out\_hook** function. This hook function is invoked after inserting the IP options from the upper protocol layers, constructing the complete IP header, and locating a route to the destination IP address. The **ip\_fltr\_out\_hook** function may modify the outgoing IP packet (*pkt*), but the interface and route have already been assigned and may not be changed. The return value from the **ip\_fltr\_out\_hook** function indicates whether *pkt* should be transmitted or dropped. See the Expected Return Values section below. If *pkt* is not dropped (**FIREWALL\_OK**), its source address is verified to be local and, if *pkt* is to be broadcast, the ability to broadcast is confirmed. Thereafter, *pkt* is enqueued on the interfaces (*ifp*) output queue. If *pkt* is dropped (**FIREWALL\_NOTOK**), it is not transmitted and **EACCES** is returned to the process.

The **inbound\_fw** and **outbound\_fw** firewall hooks allow kernel extensions to get control of packets at the place where IP receives them. If **inbound\_fw** is set, **ipintr\_noqueue**, the IP input routine, calls **inbound\_fw** and then exits. If not, **ipintr\_noqueue** calls **ipintr\_noqueue\_post\_fw** and then exits. If the **inbound\_fw** hook routine wishes to pass the packet into IP, it can call **ipintr\_noqueue\_post\_fw**. **inbound\_fw** may copy its *args* parameter by calling **inbound\_fw\_save\_args**, and may free its copy of its *args* parameter by calling **inbound\_fw\_free\_args**.

Similarly, **ip\_output** calls **outbound\_fw** if it is set, and calls **ip\_output\_post\_fw** if not. The **outbound\_fw** hook can call **ip\_output\_post\_fw** if it wants to send a packet. **outbound\_fw** may copy its *args* parameter by calling **outbound\_fw\_save\_args**, and later free its copy of its *args* parameter by calling **outbound\_fw\_free\_args**.

## Flags

Item	Description
<b>IP_FORWARDING</b>	Indicates that most of the IP headers exist.
<b>IP_RAWOUTPUT</b>	Indicates that the raw IP header exists.
<b>IP_MULTICAST_OPTS</b>	Indicates that multicast options are present.
<b>IP_ROUTETOIF</b>	Contains bypass routing tables.
<b>IP_ALLOWBROADCAST</b>	Provides capability to send broadcast packets.
<b>IP_BROADCASTOPTS</b>	Contains broadcast options inside.
<b>IP_PMTUOPTS</b>	Provides PMTU discovery options.
<b>IP_GROUP_ROUTING</b>	Contains group routing gidlist.

## Expected Return Values

Item	Description
<b>FIREWALL_OK</b>	Indicates that <i>pkt</i> is acceptable for any of the filtering functions. It will be delivered, forwarded, or transmitted as appropriate.
<b>FIREWALL_NOTOK</b>	Indicates that <i>pkt</i> should be dropped. It will not be received ( <b>ip_fltr_in_hook</b> , <b>ipsec_decap_hook</b> ) or transmitted ( <b>ip_fltr_out_hook</b> ).
<b>FIREWALL_OK_NOTSEC</b>	Indicates a return value only valid for the <b>ipsec_decap_hook</b> function. This indicates that <i>pkt</i> is acceptable according to the filtering rules, but is not encapsulated; <i>pkt</i> will be processed by the transport layer rather than processed as a decapsulated IP packet.

### Related information:

Network Kernel Services

## **i\_reset Kernel Service**

### **Purpose**

Resets a bus interrupt level.

### **Syntax**

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/intr.h>
```

```
void i_reset ( handler)
struct intr *handler;
```

### **Parameter**

Item	Description
<i>handler</i>	Specifies the address of an interrupt handler structure passed to the <i>i_init</i> service.

### **Description**

The *i\_reset* service resets the bus interrupt specified by the *handler* parameter. A device interrupt handler calls the *i\_reset* service after resetting the interrupt at the device on the bus. See *i\_init* kernel service for a brief description of interrupt handlers.

### **Execution Environment**

The *i\_reset* kernel service can be called from either the process or interrupt environment.

### **Return Values**

The *i\_reset* service has no return values.

#### **Related reference:**

“*i\_init* Kernel Service” on page 217

#### **Related information:**

Understanding Interrupts  
I/O Kernel Services

## **i\_sched Kernel Service**

### **Purpose**

Schedules off-level processing.

### **Syntax**

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/intr.h>
```

```
void i_sched ( handler)
struct intr *handler;
```

### **Parameter**

Item	Description
<i>handler</i>	Specifies the address of the pinned interrupt handler structure.

## Description

The **i\_sched** service allows device drivers to schedule some of their work to be processed at a less-favored interrupt priority. This capability allows interrupt handlers to run as quickly as possible, avoiding interrupt-processing delays and overrun conditions. See the **i\_init** kernel service for a brief description of interrupt handlers.

Processing can be scheduled off-level in the following situations:

- The interrupt handler routine for a device driver must perform time-consuming processing.
- This work does not need to be performed immediately.

**Attention:** The caller cannot alter any fields in the **intr** structure from the time the **i\_sched** service is called until the kernel calls the off-level routine. The structure must also stay pinned. Otherwise, the system may crash.

The interrupt handler structure pointed to by the *handler* parameter describes an off-level interrupt handler. The caller of the **i\_sched** service must set up all fields in the **intr** structure. The **INIT\_OFFL***n* macros in the `/usr/include/sys/intr.h` file can be used to initialize the *handler* parameter. The *n* value represents the priority class that the off-level handler should run at. Currently, classes from 0 to 3 are defined.

Use of the **i\_sched** service has two additional restrictions:

First, the **i\_sched** service will not re-register an **intr** structure that is already registered for off-level handling. Since **i\_sched** has no return value, the service will simply return normally without registering the specified structure if it was already registered but not yet executed. The kernel removes the **intr** structure from the registration list immediately prior to calling the off-level handler specified in the structure. It is therefore possible for the off-level handler to use the structure again to register another off-level request.

Care must be taken when scheduling off-level requests from a second-level interrupt handler (SLIH). If the off-level request is already registered but has not yet executed, a second registration will be ignored. If the off-level handler is currently executing, or has already run, a new request will be registered. Users of this service should be aware of these timing considerations and program accordingly.

Second, the kernel uses the `flags` field in the specified **intr** structure to determine if this structure is already registered. This field should be initialized once before the first call to the **i\_sched** service and should remain unmodified for future calls to the **i\_sched** service.

**Note:** Off-level interrupt handler path length should not exceed 5,000 instructions. If it does exceed this number, real-time support is adversely affected.

## Execution Environment

The **i\_sched** kernel service can be called from either the process or interrupt environment.

## Return Values

The **i\_sched** service has no return values.

### Related reference:

“i\_init Kernel Service” on page 217

**Related information:**

Understanding Interrupts  
I/O Kernel Services

**i\_unmask Kernel Service  
Purpose**

Enables a bus interrupt level.

**Syntax**

```
#include <sys/types.h>  
#include <sys/errno.h>  
#include <sys/intr.h>
```

```
void i_unmask ( handler)  
struct intr *handler;
```

**Parameter**

Item	Description
<i>handler</i>	Specifies the address of the interrupt handler structure that was passed to the <i>i_init</i> service.

**Description**

The *i\_unmask* service enables the bus interrupt level specified by the *handler* parameter.

**Execution Environment**

The *i\_unmask* kernel service can be called from either the process or interrupt environment.

**Return Values**

The *i\_unmask* service has no return values.

**Related reference:**

“*i\_init* Kernel Service” on page 217  
“*i\_mask* Kernel Service” on page 219

**Related information:**

Understanding Interrupts

**IS64U Kernel Service  
Purpose**

Determines if the current user-address space is 64-bit or not.

**Syntax**

```
#include <sys/types.h> #include <sys/user.h> int IS64U
```

**Description**

The *IS64U* kernel service returns 1 if the current user-address space is 64-bit. It returns 0 otherwise.

## Execution Environment

The IS64U kernel service can be called from a process or interrupt handler environment. In either case, it will operate only on the current user-address space.

## Return Values

Item	Description
0	The current user-address space is 32-bits.
1	The current user-address space is 64-bits.

### Related reference:

“as\_att64 Kernel Service” on page 11

### Related information:

Memory Kernel Services

Understanding Virtual Memory Manager Interfaces

## k

The following kernel services begin with the with the letter k.

### k\_cpuextintr\_ctl Kernel Service

#### Purpose

Performs CPU external interrupt control related operations.

#### Syntax

```
#include <sys/intr.h>
```

```
kerno_t k_cpuextintr_ctl (command , cpuset , flags)  
extintctl_t command;  
rsethandle_t cpuset;  
unit flags;
```

#### Description

This kernel services provides means of enabling, disabling, and querying the external interrupt state on the CPUs described by the CPU resource set. Enabling or disabling an CPU external interrupt could affect the external interrupt delivery to the CPU. Normally, on multiple CPU system, external interrupts can be delivered to any running CPU, and the distribution among the CPUs is determined by a predefined method. Any external interrupt can only be delivered to a CPU if its interrupt priority is more favored than the current external interrupt priority of the CPU. When external interrupts are disabled via this interface, any external interrupt priority less favored than **INTMAX** will be blocked until interrupts are enabled again. This kernel service is applicable only on selective hardware types.

**Note:** Since this kernel service change the way that interrupts are delivered, system performance may be affected. This service guarantees at least one online CPU will have external interrupts enabled for all device interrupts. Any DLPAR CPU removal can fail if the operation breaks this guarantee. On an I/O bound system, one CPU may not be enough to handle all of the external interrupts received by the partition. Performance may suffer when there are not enough CPUs enabled to handle external interrupts.

#### Parameters



Item	Description
<i>command</i>	<p>Specifies the operation to the CPU specified by the CPU resource set. One of the following values defined in <code>&lt;sys/intr.h&gt;</code> can be used:</p> <p>The following commands are supported:</p> <ul style="list-style-type: none"> <li>• <code>EXTINTDISABLE</code>: Disable external interrupts on the CPUs specified by the CPU resource set.</li> <li>• <code>EXTINTENABLE</code>: Enable external interrupts on the CPUs specified by the CPU resource set</li> <li>• <code>QUERYEXTINTDISABLE</code>: Return a CPU resource set containing the CPUs that have external interrupts disabled.</li> <li>• <code>QUERYEXTINTENABLE</code>: Return a CPU resource set containing the CPUs that have externals interrupt enabled.</li> </ul>
<i>cpuset</i>	<p>Reference to a CPU resource set. Upon successful return from this kernel service, the CPUs that have the external interrupt control operation done will be set in the CPU resource set.</p> <p>The CPUs specified by this <i>cpuset</i> parameter are logic CPU ids.</p>
<i>flags</i>	<p>Always set to 0 or <code>EINVAL_INTR_DIS_BAD_FLAGS</code> will be returned.</p>

## Security

The caller must have root authority with `CAP_NUMA_ATTACH` capability or `PV_KER_CONF` privilege in the RBAC environment.

## Execution Environment

The `k_cpuextintr_ctl` kernel service can be called from process environment only.

## Return Values

Upon successful completion, the `k_cpuextintr_ctl` kernel service returns a 0. If unsuccessful, one of the following `kerrno` value is returned.

Item	Description
<code>kerrno</code>	Description
<code>EINVAL_EXTINTR_BAD_COMMAND</code>	The command value is not valid.
<code>EINVAL_EXTINTR_BAD_FLAGS</code>	The flags value is unknown.
<code>EINVAL_EXTINTR_BAD_CPUSSET</code>	The <i>cpuset</i> references NULL.
<code>EINVAL_EXTINTR_NO_RSET</code>	The <i>cpuset</i> is empty.
<code>ENOTSUP_EXTINTR_CALLER</code>	The kernel service is called from the interrupt environment.
<code>ENOSYS_EXTINTR_PLATFORM</code>	This function is not implemented on the platform.
<code>EPERM_EXTINTR_OPER</code>	The caller does not have enough privilege to perform the requested operation.

**Note:** A return value of success does not necessarily indicate that external interrupts have been enabled or disabled on all of the specified CPUs. For example, if a CPU is not online, then the enable or disable operation will not be performed on that CPU. The caller should check the returned *cpuset* to see which CPUs have this operation successfully done. The `k_cpuextintr_ctl` kernel service will not block DR CPU add/remove operation during the whole period of system call.

## `kcap_is_set` and `kcap_is_set_cr` Kernel Service Purpose

Determines if the given capability is present in an effective capability set.

## Syntax

```
kcap_is_set (capability)  
cap_value_t capability;
```

```
kcap_is_set_cr (capability, cred)  
cap_value_t capability;  
struct ucred *cred;
```

## Parameters

Item	Description
<i>capability</i>	Specifies the capability to be examined. Must be one of the capabilities named in the <b>sys/capabilities.h</b> header file.
<i>cred</i>	Pointer to the credentials to be examined.

## Description

The **kcap\_is\_set** subroutine determines if the given capability is present in the current process' effective capability set. The **kcap\_is\_set\_cr** subroutine determines if the given capability is present in the effective capability set of the credentials structure referenced by the *cred* parameter. The *cred* parameter must be a valid referenced credentials structure.

## Return Values

The **kcap\_is\_set** and **kcap\_is\_set\_cr** subroutines return 1 if the capability is present. Otherwise, they return 0.

### Related information:

Security Kernel Services

## **kcid\_curproc** Kernel Service Purpose

Returns the current workload partition ID associated with the calling process.

## Syntax

```
#include <sys/wparid.h>
```

```
cid_t kcid_curproc ( )
```

## Description

The **kcid\_curproc** kernel service returns the workload partition ID associated with the calling process. You can use this service to determine whether the requesting process is operating within a workload partition (WPAR).

## Execution Environment

The **kcid\_curproc** kernel service can be called from the process environment only.

## Return Values

If the **kcid\_curproc** kernel service is successful, it returns the workload partition ID associated with the calling process. If the calling process is not operating within a WPAR, the ID returned is equivalent to the **WPAR\_GLOBAL** definition found in the **wparid.h** header file.

### Related reference:

“kwpar\_r2vmap\_pid Kernel Service” on page 319

“kwpar\_v2rmap\_pid Kernel Service” on page 327

## **kcred\_genpagvalue Kernel Service**

### **Purpose**

Generates a system-wide unique PAG value for a given PAG type.

### **Syntax**

```
int kcred_genpagvalue(crp, pag_type, pag_value, pag_flags);
cred_t *crp;
int pag_type;
uint64_t * pag_value;
int pag_flags;
```

### **Description**

The **kcred\_genpagvalue** kernel service generates a new PAG value for a given PAG type. It is essential that for this function to succeed the PAG type must have been previously registered with the operating system using the **kcred\_setpagname** kernel service. The scope of the **kcred\_genpagvalue** kernel service is limited to maintaining information about the last generated PAG number and accordingly generating a new number. This service optionally stores the PAG value in the **cred** structure. It does not monitor the PAG values stored in the **cred** structure by other means.

The caller must convert a PAG name to a PAG type using the **kcred\_getpagid** kernel service prior to invoking the **kcred\_genpagvalue** kernel service.

The *pag\_flags* parameter with the **PAG\_SET\_VALUE** value set causes the generated value to be atomically stored in the process's credentials.

The PAG value returned is of size 64 bits. The number of significant bits is determined by the requested PAG type. 32-bit PAGs have 32 significant bits. 64-bit PAGs have 62 significant bits.

### **Parameters**

<b>Item</b>	<b>Description</b>
<i>pag_type</i>	The <i>pag_type</i> parameter is the ID value associated with a PAG name.
<i>pag_value</i>	This pointer points to a buffer where the OS will return the newly generated PAG value.
<i>pag_flags</i>	This parameter must be 0 or the value <b>PAG_SET_VALUE</b> .

### **Return Values**

A value of 0 is returned upon successful completion. A negative value is returned if unsuccessful.

### **Error Codes**

<b>Item</b>	<b>Description</b>
EINVAL	The PAG value cannot be generated because the named PAG type does not exist as part of the table.
EPERM	The named PAG type is a 32-bit PAG and the caller does not have the SET_PROC_DAC privilege.

### **Related reference:**

“\_\_pag\_getid System Call” on page 399

“kcred\_getpagid Kernel Service” on page 246

### **Related information:**

genpagvalue Subroutine

## **kcred\_getcap Kernel Service**

### **Purpose**

Copies a capability vector from a credentials structure.

### **Syntax**

```
#include <sys/capabilities.h>
```

```
#include <sys/cred.h>
```

```
int kcred_getcap ( crp, cap )
struct ucred * cr;
struct __cap_t * cap;
```

### **Parameters**

Item	Description
<i>crp</i>	Pointer to a credentials structure
<i>cap</i>	Capabilities set

### **Description**

The **kcred\_getcap** kernel service copies the capability set from the credentials structure referenced by *crp* into *cap*. *crp* must be a valid, referenced credentials structure.

### **Execution Environment**

The **kcred\_getcap** kernel service can be called from the process environment only.

### **Return Values**

Item	Description
0	Success.
-1	An error has occurred.

### **Related information:**

Security Kernel Services

## **kcred\_getgroups Kernel Service**

### **Purpose**

Copies the concurrent group set from a credentials structure.

### **Syntax**

```
#include <sys/cred.h>
```

```
int kcred_getgroups ( crp, ngroups, groups )
struct ucred * cr;
int ngroups;
gid_t * groups;
```

### **Parameters**

Item	Description
<i>crp</i>	Pointer to a credentials structure
<i>ngroups</i>	Size of the array of group ID values
<i>groups</i>	Array of group ID values

## Description

The `kcred_getgroups` kernel service returns up to *ngroups* concurrent group set members from the credentials structure pointed to by *crp*. *crp* must be a valid referenced credentials structure.

## Execution Environment

The `kcred_getgroups` kernel service can be called from the process environment only.

## Return Values

Item	Description
$\geq 0$	The number of concurrent groups copied to groups.
-1	An error has occurred.

### Related information:

Security Kernel Services

## `kcred_getpag` or `kcred_getpag64` Kernel Service Purpose

Copies a process authentication group (PAG) ID from a credentials structure.

## Syntax

```
#include <sys/cred.h>
```

```
int kcred_getpag ( crp, which, pag )
struct ucred * cr;
int which;
int * pag;
```

```
int kcred_getpag64 ( crp, which, pag )
struct ucred * cr;
int which;
uint64 * pag;
```

## Parameters

Item	Description
<i>crp</i>	Pointer to a credentials structure
<i>which</i>	PAG ID to get
<i>pag</i>	Process authentication group

## Description

The `kcred_getpag` or `kcred_getpag64` kernel service copies the requested PAG from the credentials structure referenced by *crp* into *pag*. The value of *which* must be a defined PAG ID. The PAG ID for the *Distributed Computing Environment* (DCE) is 0. *crp* must be a valid, referenced credentials structure.

## Execution Environment

The `kcred_getpag` or `kcred_getpag64` kernel service can be called from the process environment only.

## Return Values

Upon successful completion, a value of 0 is returned. Otherwise, a value of -1 is returned, and the **errno** global variable is set to indicate the error.

## Error Codes

The `kcred_getpag` kernel service fails if the following condition is true:

Item	Description
-EOVERFLOW	PAG value is 64-bit (should be using <code>kcred_getpag64</code> )

### Related information:

Security Kernel Services

## `kcred_getpagid` Kernel Service Purpose

Returns the PAG identifier for a PAG name.

### Syntax

```
int kcred_getpagid (name)
char *name;
```

### Description

Given a PAG type name, the `kcred_getpagid` subroutine returns the PAG identifier for that PAG name.

### Parameters

Item	Description
<i>name</i>	A pointer to the name of the PAG type whose integer PAG identifier is to be returned.

## Return Values

A return value greater than or equal to 0 is the PAG identifier. A value less than 0 indicates an error.

## Error Codes

Item	Description
ENOENT	The <i>name</i> parameter doesn't refer to an existing PAG entry.

### Related reference:

“`__pag_getid` System Call” on page 399

“`__pag_getvalue` System Call” on page 400

“`kcred_getpagname` Kernel Service” on page 247

## `kcred_getpaginfo` Kernel Service Purpose

Returns a Process Authentication Group (PAG) flags for a given PAG type.

## Syntax

```
#include <sys/cred.h>
```

```
int kcred_getpaginfo ( type, infop, infosz )
int type;
struct paginfo * infop
int infosz;
```

## Parameters

Item	Description
<i>type</i>	PAG for which the flags are returned
<i>infop</i>	Pointer to PAG info structure
<i>infosz</i>	Size of <b>paginfo</b> structure

## Description

The **kcred\_getpaginfo** kernel service retrieves the flags for the specific PAG type and stores them in a PAG info structure. The value of *type* must be a defined PAG ID. The PAG ID for the Distributed Computing Environment (DCE) is 0. The *infop* parameter must be a valid, referenced PAG info structure of the size specified by *infosz*.

## Execution Environment

The **kcred\_getpaginfo** kernel service can be called from the process environment only.

## Return Values

A value of 0 is returned upon successful completion. Upon failure, a -1 is returned and **errno** is set to a value that explains the error.

### Related information:

Security Kernel Services

## **kcred\_getpagname** Kernel Service Purpose

Retrieves the name of a PAG.

## Syntax

```
int kcred_getpagname (type, buf, size)
int type;
char *buf;
int size;
```

## Description

The **kcred\_getpagname** kernel service retrieves the name of a PAG type given its integer value.

## Parameters

Item	Description
<i>type</i>	The integer valued identifier representing the PAG type.
<i>buf</i>	A <code>char *</code> to where the PAG name is copied.
<i>size</i>	An <code>int</code> that specifies the size of <i>buf</i> in bytes. The size of the buffer must be <code>PAG_NAME_LENGTH_MAX+1</code> .

## Return Values

If successful, a 0 is returned. If unsuccessful, an error code value less than 0 is returned. The PAG name associated with *type* is copied into the caller-supplied buffer *buf*.

## Error Codes

Item	Description
EINVAL	The value of <i>id</i> is less than 0 or greater than the maximum PAG identifier.
ENOENT	There is no PAG associated with <i>id</i> .
ENOSPC	The <i>size</i> parameter is insufficient to hold the PAG name.

### Related reference:

“`__pag_getid` System Call” on page 399

“`__pag_getname` System Call” on page 400

“`kcred_setpagname` Kernel Service” on page 252

## `kcred_getppriv` Kernel Service Purpose

Copies a privilege vector from a credentials structure.

### Syntax

```
#include <sys/priv.h>
#include <sys/cred.h>
```

```
int kcred_getppriv (crp, which, privset)
struct ucred *crp;
int which;
privg_t privset;
```

### Parameters

Item	Description
<i>crp</i>	Points to a credentials structure.
<i>which</i>	Specifies the privilege set to get.
<i>privset</i>	Specifies the privilege set.

### Description

The `kcred_getppriv` kernel service returns a single privilege set from the credentials structure referenced by the *crp* parameter. The *which* parameter is one of the values of `PRIV_EFFECTIVE`, `PRIV_MAXIMUM`, `PRIV_INHERITED`, `PRIV_LIMITING`, and `PRIV_USED`. The corresponding privilege set is copied to the *privset* parameter. The *crp* parameter must be a valid, referenced credentials structure.

### Execution Environment

The `kcred_getppriv` kernel service can be called from the process environment only.



## Return Values

Item	Description
0	Success.
-1	An error has occurred.

### Related information:

Security Kernel Services

## kcred\_getpriv Kernel Service Purpose

Copies a privilege vector from a credentials structure.

### Syntax

```
#include <sys/priv.h>
```

```
#include <sys/cred.h>
```

```
int kcred_getpriv ( crp, which, priv )  
struct ucred * cr;  
int which;  
priv_t * priv;
```

### Parameters

Item	Description
<i>crp</i>	Pointer to a credentials structure
<i>which</i>	Privilege set to get
<i>priv</i>	Privilege set

### Description

The `kcred_getpriv` kernel service returns a single privilege set from the credentials structure referenced by *crp*. The *which* parameter is one of `PRIV_BEQUEATH`, `PRIV_EFFECTIVE`, `PRIV_INHERITED`, or `PRIV_MAXIMUM`. The corresponding privilege set will be copied to *priv*. *rp* must be a valid, referenced credentials structure.

### Execution Environment

The `kcred_getpriv` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Success. to priv.
-1	An error has occurred.

### Related information:

Security Kernel Services

## kcred\_setcap Kernel Service Purpose

Copies a capabilities set into a credentials structure.

## Syntax

```
#include <sys/capabilities.h>
```

```
#include <sys/cred.h>
```

```
void kcred_setcap ( crp, cap )  
struct ucred * cr;  
struct __cap_t * cap;
```

## Parameters

Item	Description
<i>crp</i>	Pointer to a credentials structure
<i>cap</i>	Capabilities set

## Description

The **kcred\_setcap** kernel service initializes the capability set in the credentials structure referenced by *crp* with *cap*. *rp* must be a valid, referenced credentials structure and must not be the current credentials of any process.

## Execution Environment

The **kcred\_setcap** kernel service can be called from the process environment only.

## Return Values

The **kcred\_setcap** kernel service has no return values.

### Related information:

Security Kernel Services

## kcred\_setgroups Kernel Service Purpose

Copies a concurrent group set into a credentials structure.

## Syntax

```
#include <sys/cred.h>
```

```
int kcred_setgroups ( crp, ngroups, groups )  
struct ucred * cr;  
int ngroups;  
gid_t * groups;
```

## Parameters

Item	Description
<i>crp</i>	Pointer to a credentials structure
<i>ngroups</i>	Size of the array of group ID values
<i>groups</i>	Array of group ID values

## Description

The **kcred\_setgroups** kernel service copies *ngroups* concurrent group set members into the credentials structure pointed to by *crp*. *crp* must be a valid, referenced credentials structure and must not be the current credentials of any process.

## Execution Environment

The `kcred_setgroups` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	The concurrent group set has been copied successfully.
-1	An error has occurred.

### Related information:

Security Kernel Services

## `kcred_setpag` or `kcred_setpag64` Kernel Service Purpose

Copies a process authentication group ID into a credentials structure.

## Syntax

```
#include <sys/cred.h>
```

```
int kcred_setpag ( crp, which, pag )
struct ucred * cr;
int which;
int pag;
```

```
int kcred_setpag64 ( crp, which, pag )
struct ucred * cr;
int which;
uint64 * pag;
```

## Parameters

Item	Description
<i>crp</i>	Pointer to a credentials structure
<i>which</i>	PAG ID to set
<i>pag</i>	Process authentication group

## Description

The `kcred_setpag` or `kcred_setpag64` kernel service initializes the specified PAG in the credentials structure referenced by *crp* with *pag*. The value of *which* must be a defined PAG ID. The PAG ID for the *Distributed Computing Environment* (DCE) is 0. *Crp* must be a valid, referenced credentials structure. *crp* may be a reference to the current credentials of a process.

## Execution Environment

The `kcred_setpag` or `kcred_setpag64` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Success.
-1	An error has occurred.

#### Related information:

Security Kernel Services

## kcred\_setpagname Kernel Service Purpose

Copies a process authentication group ID into a credentials structure.

### Syntax

```
int kcred_setpagname (name, flags, func)
char *name;
int flags;
```

### Description

The **kcred\_setpagname** kernel service registers the name of a PAG and returns the PAG type identifier. If the PAG name has already been registered, the previously returned PAG type identifier is returned if the *flags* and *func* parameters match their earlier values.

### Parameters

Item	Description
<i>name</i>	The <i>name</i> parameter is a 1 to 4 character, NULL-terminated name for the PAG type. Typical values might include "afs", "dfs", "pki" and "krb5."
<i>flags</i>	The <i>flags</i> parameter indicates if each PAG value is unique (PAG_UNIQUEVALUE) or multivalued (PAG_MULTIVALUED). A multivalued PAG type allows multiple calls to the <b>kcred_setpag</b> kernel service to be made to store multiple values for a single PAG type.
<i>func</i>	The <i>func</i> parameter is a pointer to an allocating and deallocating function. The <i>flag</i> parameter to that function is either PAGVALUE_ALLOC or PAGVALUE_FREE. The <i>value</i> parameter is the actual PAG value. The <i>func</i> parameter will be invoked by the <b>crfree</b> kernel service with a flag value of PAGVALUE_FREE on the last free value of a credential. Whenever a credentials structure is initialized with new PAG values, <i>func</i> will be invoked by that function with a value of PAGVALUE_ALLOC. This parameter may be ignored and an error returned if the value of <i>func</i> is non-NULL.

### Return Values

A value of 0 or greater is returned upon successful completion. This value is the PAG type identifier which is used with other kernel services, such as the **kcred\_getpag** and **kcred\_setpag** subroutines . A negative value is returned if unsuccessful.

### Error Codes

Item	Description
ENOSPC	The PAG table is full.
EEXISTS	The named PAG type already exists in the table and the <i>flags</i> and <i>func</i> parameters do not match their earlier values.
EINVAL	The <i>flags</i> parameter is an invalid value.

#### Related reference:

"\_\_pag\_setname System Call" on page 401

"\_\_pag\_setvalue System Call" on page 402

"kcred\_getpagname Kernel Service" on page 247

## kcred\_setppriv Kernel Service Purpose

Copies a privilege vector into a credentials structure.

### Syntax

```
#include <sys/priv.h>
#include <sys/cred.h>
```

```
int kcred_setppriv ( crp, which, privset )
struct ucred *crp;
int which;
privg_t privset;
```

### Parameters

Item	Description
<i>crp</i>	Points to a credentials structure.
<i>which</i>	Specifies the privilege set to set.
<i>privset</i>	Specifies the privilege set.

### Description

The **kcred\_setppriv** kernel service sets one or more single privilege sets in the credentials structure referenced by the *crp* parameter. The *which* parameter is the bitwise OR of one or more values of **PRIV\_EFFECTIVE**, **PRIV\_MAXIMUM**, **PRIV\_INHERITED**, **PRIV\_LIMITING**, and **PRIV\_USED**. The *privset* parameter initializes the corresponding privilege sets. The *crp* parameter must be a valid, referenced credentials structure and cannot be the current credentials of any process.

### Execution Environment

The **kcred\_setppriv** kernel service can be called from the process environment only.

### Return Values

Item	Description
0	Success.
-1	An error has occurred.

### Related information:

Security Kernel Services

## kcred\_setpriv Kernel Service Purpose

Copies a privilege vector into a credentials structure.

### Syntax

```
#include <sys/priv.h>
#include <sys/cred.h>
```

```
int kcred_setpriv ( crp, which, priv )
struct ucred * cr;
int which;
priv_t * priv;
```

## Parameters

Item	Description
<i>crp</i>	Pointer to a credentials structure
<i>which</i>	Privilege set to set
<i>priv</i>	Privilege set

## Description

The `kcred_setpriv` kernel service sets one or more single privilege sets in the credentials structure referenced by *crp*. The *which* parameter is one or more bit-wise ored values of `PRIV_BEQUEATH`, `PRIV_EFFECTIVE`, `PRIV_INHERITED`, and `PRIV_MAXIMUM`. The corresponding privilege sets are initialized from *priv*. *crp* must be a valid, referenced credentials structure and must not be the current credentials of any process.

## Execution Environment

The `kcred_setpriv` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Success. to priv.
-1	An error has occurred.

### Related information:

Security Kernel Services

## kern\_soaccept Kernel Service Purpose

Accepts the first queued connection by assigning it to the new socket.

## Syntax

```
#include <sys/kern_socket.h>
int kern_soaccept( ksocket_t so,
ksocket_t *aso,
struct mbuf **name,
int nonblock )
```

## Parameters

Item	Description
<i>so</i>	The socket that is used in the <code>kern_solisten()</code> Kernel Service.
<i>aso</i>	The new socket for the accepted connection. The caller must pass in the address of the <code>ksocket_t</code> .
<i>name</i>	A <code>struct sockadr</code> address is returned in a <code>mbuf</code> buffer whose address is stored in the <code>*name</code> parameter. The caller should pass in the address of the <code>struct mbuf *</code> structure. The caller sets the <code>mbuf</code> buffer free after the function returns successfully.
<i>nonblock</i>	A flag to specify if this call should be nonblocking. The value of 1 is for nonblocking and 0 is for blocking.

## Description

The `kern_soaccept` kernel service accepts the first queued connection by assigning it to the new socket.

## Execution Environment

The `kern_socket` kernel service can be called from the process environment.

### Examples

```
struct mbuf *name = NULL;
ksocket_t so;
ksocket_t aso;
struct sockaddr_in laddr;
int rc;
rc = kern_socket(AF_INET, &so, SOCK_STREAM, IPPROTO_TCP);
if (rc != 0 )
{
    return(-1);
}
bzero(&laddr, sizeof(struct sockaddr_in));
laddr.sin_family = AF_INET;
laddr.sin_port = 12345;
laddr.sin_len = sizeof(struct sockaddr_in);
laddr.sin_addr.s_addr = inet_addr("9.3.108.208");
rc = kern_socketbind(so, (struct sockaddr *)&laddr);
if (rc != 0 )
{
    return(-1);
}
rc = kern_socketlisten(so, 5);
if (rc != 0 )
{
    return(-1);
}
rc = kern_socketaccept(so, &aso, &name, 0);
if (rc != 0 )
{
    return(-1);
}
m_freem(name); /* Caller needs to free the mbuf after kern_socketaccept */
```

### Return Values

Item	Description
0	Upon Success
>0	Error

The non-zero return value is the error number that is defined in the `/usr/include/sys/errno.h` file.

#### Related reference:

“`kern_socket` Kernel Service” on page 258

“`kern_socketreceive` Kernel Service” on page 261

“`kern_socketsend` Kernel Service” on page 264

## `kern_socketbind` Kernel Service

### Purpose

Associates the local network address to the socket.

### Syntax

```
#include <sys/kern_socket.h>
int kern_socketbind( ksocket_t so, struct sockaddr *laddr )
```

## Parameters

Item	Description
<i>so</i>	The socket that was created by the <code>kern_socreate()</code> system call.
<i>laddr</i>	Local address to be bound.

## Description

The `kern_sobind` kernel service binds a local address to the socket.

## Execution Environment

The `kern_sobind` kernel service can be called from the process environment.

## Examples

```
ksocket_t  so;
struct sockaddr_in laddr;
int        rc;
rc = kern_socreate(AF_INET, &so, SOCK_STREAM, IPPROTO_TCP);
if (rc != 0 )
{
    return(-1);
}
bzero(&laddr, sizeof(struct sockaddr_in));
laddr.sin_family = AF_INET;
laddr.sin_port = 12345;
laddr.sin_len = sizeof(struct sockaddr_in);
laddr.sin_addr.s_addr = inet_addr("9.3.108.208");
rc = kern_sobind(so, (struct sockaddr *) &laddr);
if (rc != 0 )
{
    return(-1);
}
```

## Return Values

Item	Description
0	Upon Success
>0	Error

The nonzero return value is the error number that is defined in the `/usr/include/sys/errno.h` file.

### Related reference:

“`kern_socreate` Kernel Service” on page 258

“`kern_solisten` Kernel Service” on page 260

## kern\_soclose Kernel Service

### Purpose

Aborts any connections and releases the data in the socket.

### Syntax

```
#include <sys/kern_socket.h>
int kern_soclose( ksocket_t so )
```

## Parameters



Item	Description
<code>so</code>	The socket on which the close will be issued.

## Description

The `kern_soclose` kernel service aborts any connection and releases the data in the socket.

## Execution Environment

The `kern_soclose` kernel service can be called from the process environment.

## Examples

```
ksocket_t so;
int rc;
rc = kern_socreate(AF_INET, &so, SOCK_STREAM, IPPROTO_TCP);
if (rc != 0 )
{
    return(-1);
}
/* Socket is in use */
...
kern_soclose(so);
```

## Return Values

Item	Description
0	Upon Success
>0	Error

The nonzero return value is the error number that is defined in the `/usr/include/sys/errno.h` file.

### Related reference:

“`kern_socreate` Kernel Service” on page 258

## `kern_soconnect` Kernel Service

### Purpose

Establishes a connection to a foreign address.

### Syntax

```
#include <sys/kern_socket.h>
int kern_soconnect( ksocket_t so, struct sockaddr *faddr )
```

### Parameters

Item	Description
<code>so</code>	The socket that was created by <code>socreate()</code> .
<code>faddr</code>	Foreign address to connect.

## Description

The `kern_soconnect` kernel service establishes connection with a foreign address.

## Execution Environment

The `kern_soconnect` kernel service can be called from the process environment.

## Examples

```
ksocket_t so;
struct sockaddr_in faddr;
int rc;
rc = kern_screate(AF_INET, &so, SOCK_STREAM, IPPROTO_TCP);
if (rc != 0 )
{
    return(-1);
}
bzero(&faddr, sizeof(struct sockaddr_in));
faddr.sin_family = AF_INET;
faddr.sin_port = 23456;
faddr.sin_len = sizeof(struct sockaddr_in);
faddr.sin_addr.s_addr = inet_addr("9.3.108.210");
rc = kern_soconnect(so, (struct sockaddr *) &faddr);
if (rc != 0 )
{
    return(-1);
}
```

## Return Values

Item	Description
0	Upon Success
>0	Error

The nonzero return value is the error number that is defined in the `/usr/include/sys/errno.h` file.

### Related reference:

“kern\_screate Kernel Service”

“kern\_sosend Kernel Service” on page 264

“kern\_soreceive Kernel Service” on page 261

## kern\_screate Kernel Service

### Purpose

Used to create a socket of the specified address family and type. If the protocol is left unspecified (zero), then the system selects the protocol based on the address family and type.

### Syntax

```
#include <sys/kern_socket.h>
```

```
int kern_screate (int addressfamily, ksocket_t *so, int type, int protocol)
```

### Parameters

Item	Description
<i>addressfamily</i>	The address family for the newly created socket. The file <code>&lt;sys/socket.h&gt;</code> contains the definitions for the family. Currently AIX supports:  <b>AF_INET</b> Denotes the IPv4 Internet addresses.  <b>AF_INET6</b> Denotes the IPv6 Internet addresses.
<i>so</i>	The socket assigned by the <b>create(0)</b> call. The caller must pass the address of <b>ksocket_t</b> .
<i>type</i>	The requested socket type. The file <code>&lt;sys/socket.h&gt;</code> contains the definition for the socket type. Currently AIX supports <b>SOCK_STREAM</b> .
<i>protocol</i>	The file <code>&lt;netinet/in.h&gt;</code> contains the definition for the protocol. Currently AIX supports <b>IPPROTO_TCP</b>

## Description

The `kern_socreate` kernel service creates a socket based on the address family, type and protocol.

## Execution Environment

The `kern_socreate` kernel service can be called from the process environment.

## Examples

```
ksocket_t    so;
ksocket_t    so2;
kern_socreate(AF_INET, &so, SOCK_STREAM, IPPROTO_TCP);
kern_socreate(AF_INET6, &so2, SOCK_STREAM, 0);
```

## Return Values

Item	Description
0	Upon Success
>0	Error

The nonzero return value is the error number that is defined in the `/usr/include/sys/errno.h` file.

### Related reference:

“`kern_soclose` Kernel Service” on page 256

“`kern_soconnect` Kernel Service” on page 257

“`kern_soshutdown` Kernel Service” on page 266

## kern\_sogetopt Kernel Service Purpose

Obtains the option associated with the socket, either at the socket level or at the protocol level.

## Syntax

```
#include <sys/kern_socket.h>
int kern_sogetopt( ksocket_t so, int level, int optname, struct mbuf **mp )
```

## Parameters

Item	Description
<i>so</i>	The socket that will be used to retrieve the option.
<i>level</i>	The socket level (e.g. <code>SOL_SOCKET</code> ) or protocol level ( <code>IPPROTO_TCP</code> )
<i>optname</i>	The option name to retrieve. Socket options can be found in <code>&lt;sys/socket.h&gt;</code> and TCP options can be found in <code>&lt;netinet/tcp.h&gt;</code>
<i>mp</i>	The <code>mbuf</code> that will be returned with the option value. The <code>mp-&gt;m_len</code> will be the size of the value. The caller must pass the address of the <code>struct mbuf *</code> . The caller must set the <code>mbuf</code> free after the function returns successfully.

## Description

The `kern_sogetopt` kernel service obtains the option associated with the socket, either at the socket level, or at the protocol level.

## Execution Environment

The `kern_sogetopt` kernel service can be called from the process environment.

## Examples

```
ksocket_t so;
int rc;
struct mbuf *sopt = NULL;
int tcp_nodelay = -1;
rc = kern_socreate(AF_INET, &so, SOCK_STREAM, IPPROTO_TCP);
if (rc != 0)
{
    return(-1);
}
rc = sogetopt(so, IPPROTO_TCP, TCP_NODELAY, &sopt);
if (rc != 0)
{
    return(-1);
}
tcp_nodelay = *mtod(sopt, int *) ? 1 : 0;
m_free(sopt); /* Caller needs to free the mbuf after kern_sogetopt */
```

## Return Values

Item	Description
0	Upon Success
>0	Error

The nonzero return value is the error number that is defined in the `/usr/include/sys/errno.h` file.

### Related reference:

“kern\_socreate Kernel Service” on page 258

## kern\_solisten Kernel Service

### Purpose

Prepares to accept incoming connections on the socket.

### Syntax

```
#include <sys/kern_socket.h>
int kern_solisten(ksocket_t so, int backlog)
```

### Parameters

Item	Description
<i>so</i>	The socket that was created by <code>kern_socreate()</code> and used in <code>kern_sobind()</code>
<i>backlog</i>	Limit the number of connection requests that can be queued on this socket. The maximum value that can be passed to this parameter equals the minimum number of user backlog number and the network option <code>somaxconn</code> value.

### Description

The `kern_solisten` kernel service prepares to accept incoming connection on the socket.

### Execution Environment

The `kern_solisten` kernel service can be called from the process environment.

### Examples

```
struct mbuf *name = NULL;
ksocket_t so;
struct sockaddr_in laddr;
int rc;
```

```

rc = kern_screate(AF_INET, &so, SOCK_STREAM, IPPROTO_TCP);
if (rc != 0 )
{
    return(-1);
}
bzero(&laddr, sizeof(struct sockaddr_in));
laddr.sin_family = AF_INET;
laddr.sin_port = 12345;
laddr.sin_len = sizeof(struct sockaddr_in);
laddr.sin_addr.s_addr = inet_addr("9.3.108.208");
rc = kern_sobind(so, (struct sockaddr *)&laddr);
if (rc != 0 )
{
    return(-1);
}
rc = kern_solisten(so, 5);
if (rc != 0 )
{
    return(-1);
}

```

## Return Values

Item	Description
0	Upon Success
>0	Error

The nonzero return value is the error number that is defined in the `/usr/include/sys/errno.h` file.

### Related reference:

“kern\_screate Kernel Service” on page 258

“kern\_soaccept Kernel Service” on page 254

## kern\_soreceive Kernel Service

### Purpose

The routine processes one record per call and returns the number of bytes requested.

### Syntax

```

#include <sys/kern_socket.h>
int kern_soreceive( ksocket_t so,
struct mbuf **paddr,
long len,
struct mbuf **mp,
struct mbuf **controlp,
int *flagp )

```

### Parameters

Item	Description
<i>so</i>	The socket to receive the data.
<i>paddr</i>	The foreign socket address information is returned in this pointer. Caller should pass in address of <b>struct mbuf *</b> . Caller needs to free this <b>mbuf</b> after the function returns successfully. Caller can pass in NULL if caller doesn't need foreign address information.
<i>len</i>	The length of the data to be received.
<i>mp</i>	The <b>mbuf</b> pointer so that data can be returned in an <b>mbuf</b> chain. The caller must pass in the address of <b>struct mbuf *</b> . The caller must free this <b>mbuf</b> after the function returns.
<i>controlp</i>	Pointer to an <b>mbuf</b> containing the control information. Caller should pass in address of <b>struct mbuf *</b> . Caller needs to free this <b>mbuf</b> after the function returns successfully. Caller can pass in NULL if there is no control information.

Item	Description
<i>flagp</i>	If <i>flagp</i> is not NULL, caller can pass in actual flag. The flags are defined in the <sys/socket.h> file. The <b>kern_soreceive</b> routine will use flags set in <i>flagp</i> . The caller can set the <i>flagp</i> to <b>MSG_WAITALL</b> or <b>MSG_NONBLOCK</b> . On return, it will set <i>flagp</i> to <b>MSG_TRUNC</b> , <b>MSG_OOB</b> wherever applicable.

## Description

The **kern\_soreceive** kernel service processes one record per call and returns the number of bytes that are requested. If there is data in the socket receive buffer, the **kern\_soreceive** kernel service returns up to < len> bytes as a **mbuf** chain. The actual number of bytes returned is computed by adding the **m\_len** fields of each **mbuf** in the chain. If there is no data, but the connection is still established, the **kern\_soreceive** kernel service either returns EAGAIN with \*mp set to NULL (if MSG\_NONBLOCK is set) or returns wait for data to arrive (if MSG\_NONBLOCK is not set). If the connection is closed before the call or while waiting for data, the \*mp is set to NULL and 0 is returned. Waiting may be interrupted, in which case **kern\_soreceive** returns EAGAIN, EINTR, or ERESTART and \*mp is undefined. The application might return EINTR, but calls the **kern\_soreceive** kernel service again.

## Execution Environment

The **kern\_soreceive** kernel service can be called from the process environment.

## Examples

```
ksocket_t  so;
struct mbuf *data = NULL;
struct sockaddr_in faddr;
long       len = 512;
int        flags = 0;
int        rc;
rc = kern_socreate(AF_INET, &so, SOCK_STREAM, IPPROTO_TCP);
if (rc != 0 )
{
    return(-1);
}
bzero(&faddr, sizeof(struct sockaddr_in));
faddr.sin_family = AF_INET;
faddr.sin_port = 23456;
faddr.sin_len = sizeof(struct sockaddr_in);
faddr.sin_addr.s_addr = inet_addr("9.3.108.210");
rc = kern_soconnect(so, (struct sockaddr *) &faddr);
if (rc != 0 )
{
    return(-1);
}
do
{
    rc = kern_soreceive(so, NULL, len, &data, NULL, &flags);
} while (rc == EAGAIN || rc == EINTR || rc == ERESTART);
if ((rc == 0) && data)
{
    /* process the data */
    ...
    m_freem(data); /* Caller needs to free the mbuf after kern_soreceive. */
}
else
{
    return(-1);
}
```

## Return Values

Item	Description
0	Upon Success
>0	Error

The nonzero return value is the error number that is defined in the `/usr/include/sys/errno.h` file.

#### Related reference:

“`kern_socreate` Kernel Service” on page 258

“`kern_sosend` Kernel Service” on page 264

## kern\_soreserve Kernel Service

### Purpose

The routine enforces the limit for the send and receive buffer space for a socket. It does not actually allocate memory only sets the buffer size.

### Syntax

```
#include <sys/kern_socket.h>
int kern_soreserve( ksocket_t so, uint64_t sndcc, uint64_t rcvcc )
```

### Parameters

Item	Description
<i>so</i>	The socket that will be used in reserving the space.
<i>sndcc</i>	Send buffer size in bytes.
<i>rcvcc</i>	Receive buffer size in bytes.

### Description

The `kern_soreserve` kernel service enforces the limit for the send and receive buffer space for a socket. It does not actually allocate memory. It sets the buffer size.

### Execution Environment

The `kern_soreserve` kernel service can be called from the process environment.

### Examples

```
ksocket_t so;
uint64_t sb_snd_hiwat = 2048;
uint64_t sb_rcv_hiwat = 2048;
int rc;
rc = kern_socreate(AF_INET, &so, SOCK_STREAM, IPPROTO_TCP);
if (rc != 0)
{
    return(-1);
}
rc = kern_soreserve(so, sb_snd_hiwat, sb_rcv_hiwat);
if (rc != 0)
{
    return(-1);
}
```

### Return Values

Item	Description
0	Upon Success
>0	Error

The nonzero return value is the error number that is defined in the `/usr/include/sys/errno.h` file.

#### Related reference:

“kern\_socreate Kernel Service” on page 258

## kern\_sosend Kernel Service

### Purpose

Pass data and control information to the protocol associated send routines.

### Syntax

```
#include <sys/kern_socket.h>
int kern_sosend( ksocket_t so, struct sockaddr *faddr,
struct mbuf *top,
struct mbuf *control,
int flags )
```

### Parameters

Item	Description
<i>so</i>	The socket to send data.
<i>faddr</i>	The destination address, only necessary if the socket is not connected.
<i>top</i>	The <b>mbuf</b> chain of data to be sent. Remember that the first <b>mbuf</b> must have the packet header filled out. Set the <b>top-&gt;m_pkthdr.len</b> to the total length of the data in the <b>mbuf</b> chain and the <b>m_flags</b> to <b>M_PKTHDR</b> . The caller must allocate <b>mbuf</b> memory before calling the routine.
<i>control</i>	Pointer to an <b>mbuf</b> containing the control information to be sent. The caller must allocate <b>mbuf</b> memory before calling the function if the caller wants to pass in control information.
<i>flags</i>	Flags options for this write call. Caller can set flags to <b>MSG_NONBLOCK</b> .

### Description

The **kern\_sosend** kernel service passes data and control information to the protocol associated send routines.

### Execution Environment

The **kern\_sosend** kernel service can be called from process environment.

### Examples

```
ksocket_t so;
int flags = 0;
struct sockaddr_in faddr;
struct mbuf *send_mbuf;
struct sockaddr_in faddr;
char msg[100];
int i, rc;
rc = kern_socreate(AF_INET, &so, SOCK_STREAM, IPPROTO_TCP);
if (rc != 0 )
{
return(-1);
}
bzero(&faddr, sizeof(struct sockaddr_in));
faddr.sin_family = AF_INET;
faddr.sin_port = 23456;
faddr.sin_len = sizeof(struct sockaddr_in);
```



```

faddr.sin_addr.s_addr = inet_addr("9.3.108.210");
rc = kern_soconnect(so, (struct sockaddr *) &faddr);
if (rc != 0 )
{
    return(-1);
}
send_mbuf = MGETBUF(sizeof(msg), M_DONOTWAIT); /* Caller needs to allocate mbuf memory */
if (send_mbuf == NULL)
{
    return (-1);
}
for (i=0; i < 100, i++)
{
    msg[i] = 0x2A;
}
bcopy(msg, mtod(send_mbuf, caddr_t), sizeof(msg));
send_mbuf->m_len = send_mbuf->m_pkthdr.len = sizeof(msg);
rc = kern_sosend(so, NULL, send_mbuf, 0, MSG_NONBLOCK);
if (rc != 0 )
{
    return(-1);
}

```

## Return Values

Item	Description
0	Upon Success
>0	Error

The nonzero return value is the error number that is defined in the `/usr/include/sys/errno.h` file.

### Related reference:

“kern\_socreate Kernel Service” on page 258

“kern\_soreceive Kernel Service” on page 261

## kern\_sosetopt Kernel Service

### Purpose

Sets the option associated with the socket, either at the socket level or at the protocol level.

### Syntax

```

#include <sys/kern_socket.h>
int sosetopt( ksocket_t so,
int level, int optname,
struct mbuf *mp )

```

### Parameters

Item	Description
<i>so</i>	The socket that will be used to set the option.
<i>level</i>	The socket level (e.g. <code>SOL_SOCKET</code> ) or protocol level ( <code>IPPROTO_TCP</code> )
<i>optname</i>	The option name to set. Socket options can be found in <code>&lt;sys/socket.h&gt;</code> and the TCP options can be found in <code>&lt;netinet/tcp.h&gt;</code> .
<i>mp</i>	The <b>mbuf</b> that contains the option value and will be used to modify the field specified by the option name. The <b>mp-&gt;m_len</b> should be the size of the value. The caller must allocate <b>mbuf</b> memory before calling the function.

## Description

The `kern_sosetopt` kernel service sets the option associated with the socket, either at the socket level, or at the protocol level.

## Execution Environment

The `kern_sosetopt` kernel service can be called from the process environment.

## Examples

```
ksocket_t    so;
struct mbuf *mp = NULL;
struct linger *linger;
int rc;
rc = kern_socreate(AF_INET, &so, SOCK_STREAM, IPPROTO_TCP);
if (rc != 0 )
{
    return(-1);
}
mp = m_get(M_DONTWAIT, MT_SOOPTS); /* Caller of kern_sosetopt needs to allocated mbuf memory */
if (mp == NULL)
{
    return (-1);
}
mp->m_len = sizeof(struct linger);
linger = mtod(mp, struct linger *);
linger->l_linger = 5;
linger->l_onoff = 1;
rc = kern_sosetopt(so, SOL_SOCKET, SO_LINGER, mp);
if (rc != 0 )
{
    return(-1);
}
```

## Return Values

Item	Description
0	Upon Success
>0	Error

The nonzero return value is the error number that is defined in the `/usr/include/sys/errno.h` file.

### Related reference:

“`kern_socreate` Kernel Service” on page 258

## `kern_soshutdown` Kernel Service

### Purpose

Closes the read-half, write-half or both read and write of a connection.

### Syntax

```
#include <sys/kern_socket.h>
int kern_soshutdown( ksocket_t so, int how )
```

### Parameters

Item	Description
<i>so</i>	The socket to which the shutdown will be issued.
<i>how</i>	0 read, 1 write, 2 read and write

## Description

The `kern_soshutdown` kernel service closes the read-half, write-half or both read and write of a connection.

## Execution Environment

The `kern_soshutdown` kernel service can be called from the process environment.

## Examples

```
ksocket_t    so;

/* Create the socket so */
kern_socreate(AF_INET, &so, SOCK_STREAM, IPPROTO_TCP);

/* Shutting down both the read/write */
kern_soshutdown(so, 2);
```

## Return Values

Item	Description
0	Upon Success
>0	Error

## Related reference:

“`kern_socreate` Kernel Service” on page 258

## kgethostname Kernel Service

### Purpose

Retrieves the name of the current host.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
int
kgethostname ( name, namelen)
char *name;
int *namelen;
```

### Parameters

Item	Description
<i>name</i>	Specifies the address of the buffer in which to place the host name.
<i>namelen</i>	Specifies the address of a variable in which the length of the host name will be stored. This parameter should be set to the size of the buffer before the <b>kgethostname</b> kernel service is called.

## Description

The **kgethostname** kernel service returns the standard name of the current host as set by the **sethostname** subroutine. The returned host name is null-terminated unless insufficient space is provided.

## Execution Environment

The **kgethostname** kernel service can be called from either the process or interrupt environment.

## Return Value

Item	Description
0	Indicates successful completion.

### Related information:

sethostname subroutine

Network Kernel Services

## kgetpname Kernel Service Purpose

Provides the calling process's base program name.

## Syntax

```
#include <sys/encap.h>
int kgetpname (char * Buffer, size_t *BufferSize);
```

## Description

The **kgetpname** kernel service copies the program name of the calling process into the buffer specified by *Buffer*. Including the null terminator, the service copies no more than the lesser of *\*BufferSize*, **MAXCOMLEN**, or the actual size of the program name in bytes into the buffer. If *Buffer* is NULL, or *\*BufferSize* is 0, no copy is performed. If the full program name is copied into the buffer, the total number of bytes copied is written to *\*BufferSize*. If **kgetpname** cannot copy the full program name into the buffer, the size in bytes of the full program name is written to *\*BufferSize*, and **ENAMETOOLONG** is returned.

## Execution Environment

The **kgetpname** kernel service can only be called from the process environment.

## Return Values

Item	Description
0	The full program name was successfully written to the buffer.
ENAMETOOLONG	Only part of the full program name was written to the buffer, and <b>kgetpname</b> stored the (positive) length in bytes (including the null character) of the full program name into <i>*BufferSize</i> .
EINVAL	<i>Buffer</i> is Null, <i>BufferSize</i> is NULL, or <i>*BufferSize</i> is 0.
ENOTSUP	The <b>kgetpname</b> kernel service was called from inside an interrupt context.

## **kgetrlimit64 Kernel Service**

### **Purpose**

Controls maximum system resource consumption.

### **Library**

Standard C Library (**libc.a**)

### **Syntax**

```
#include <sys/time.h>
#include <sys/resource.h>
```

```
void kgetrlimit64 (Resource1, RLP)
int Resource1;
struct rlimit64 *RLP;
```

## Parameters

Item	Description
<i>Resource1</i>	The <i>Resource1</i> parameter can be one of the following values: <b>RLIMIT_AS</b> The maximum size of a process's total available memory, in bytes. This limit is not enforced. <b>RLIMIT_CORE</b> The largest size, in bytes, of a core file that can be created. This limit is enforced by the kernel. If the value of the RLIMIT_FSIZE limit is less than the value of the RLIMIT_CORE limit, the system uses the RLIMIT_FSIZE limit value as the soft limit. <b>RLIMIT_CPU</b> The maximum amount of central processing unit (CPU) time, in seconds, to be used by each process. If a process exceeds its soft CPU limit, the kernel sends a SIGXCPU signal to the process. After the hard limit is reached, the process is killed with SIGXCPU, even if it handles, blocks, or ignores that signal. <b>RLIMIT_DATA</b> The maximum size, in bytes, of the data region for a process. This limit defines how far a program can extend its break value with the <b>sbrk</b> subroutine. This limit is enforced by the kernel. <b>RLIMIT_FSIZE</b> The largest size, in bytes, of any single file that can be created. When a process attempts to write, truncate, or clear beyond its soft RLIMIT_FSIZE limit, the operation fails with the <b>errno</b> variable set to EFBIG. If the environment variable XPG_SUS_ENV=ON is set in the user's environment before the process is issued, then the SIGXFSZ signal is also generated. <b>RLIMIT_NOFILE</b> This is a number one greater than the maximum value that the system can assign to a newly-created descriptor. <b>RLIMIT_STACK</b> The maximum size, in bytes, of the stack region for a process. This limit defines how far a program stack region can be extended. The system automatically performs stack extension. This limit is enforced by the kernel. When the stack limit is reached, the process receives a SIGSEGV signal. If this signal is not caught by a handler using the signal stack, the signal ends the process. <b>RLIMIT_RSS</b> The maximum size, in bytes, to which the resident set size of a process can grow. This limit is not enforced by the kernel. A process might exceed its soft limit size without being ended.
<i>RLP</i>	Points to the <b>rlimit64</b> structure where the requested limits are returned by the <b>kgetrlimit64</b> kernel service.

## Description

The **kgetrlimit64** kernel service returns the values of limits on system resources used by the current process and its children processes.

**Note:** The initial values returned by the **kgetrlimit64** kernel service are the ulimit values in effect when the process was started. For maxdata programs the initial soft limit for data is set to the lower of data ulimit value or a value corresponding to the number of data segments reserved for data segments.

The **rlimit64** structure specifies the hard and soft limits for a resource, as defined in the **sys/resource.h** file. The RLIM64\_INFINITY value defines an infinite value for a limit.

## Execution Environment

The **kgetrlimit64** kernel service can be called from either the process or interrupt environment.

## Return Values

The **kgetrlimit64** kernel service has no return values.

**Related information:**

getrlimit64 subroutine

## **kgetsystemcfg Kernel Service**

### **Purpose**

Displays the system configuration information.

### **Syntax**

```
#include <systemcfg.h>
uint64_t kgetsystemcfg ( int name)
```

### **Description**

Displays the system configuration information.

### **Parameters**

Item	Description
<i>name</i>	Specifies the system variable setting to be returned. Valid values for the <i>name</i> parameter are defined in the <code>systemcfg.h</code> file.

### **Return Values**

If the value specified by the *name* parameter is system-defined, the **kgetsystemcfg** kernel service returns the data that is associated with the structure member represented by the *input* parameter. Otherwise, the **kgetsystemcfg** kernel service will return `UINT64_MAX`, and the **errno** will be set.

### **Error Codes**

The **kgetsystemcfg** subroutine will fail if:

Item	Description
<code>EINVAL</code>	The value of the <i>name</i> parameter is invalid.

### **Related information:**

getsystemcfg subroutine

## **kgettickd Kernel Service**

### **Purpose**

Retrieves the current status of the systemwide time-of-day timer-adjustment values.

### **Syntax**

```
#include <sys/types.h>
int kgettickd (timed, tickd, time_adjusted)
int *timed;
int *tickd;
int *time_adjusted;
```

### **Parameters**

Item	Description
<i>timed</i>	Specifies the current amount of time adjustment in microseconds remaining to be applied to the systemwide timer.
<i>tickd</i>	Specifies the time-adjustment rate in microseconds.
<i>time_adjusted</i>	Indicates if the systemwide timer has been adjusted. A value of True indicates that the timer has been adjusted by a call to the <b>adjtime</b> or <b>settimer</b> subroutine. A value of False indicates that it has not. The use of the <b>ksettimer</b> kernel service has no effect on this flag. This flag can be changed by the <b>ksettickd</b> kernel service.

## Description

The **kgettickd** kernel service provides kernel extensions with the capability to determine if the **adjtime** or **settimer** subroutine has adjusted or changed the systemwide timer.

The **kgettickd** kernel service is typically used only by kernel extensions providing time synchronization functions. This includes coordinated network time (which is the periodic synchronization of all system clocks to a common time by a time server or set of time servers on a network), where use of the **adjtime** subroutine is insufficient.

## Execution Environment

The **kgettickd** kernel service can be called from either the process or interrupt environment.

## Return Values

The **kgettickd** service always returns a value of 0.

### Related reference:

“ksettimer Kernel Service” on page 309

### Related information:

Timer and Time-of-Day Kernel Services

Using Fine Granularity Timer Services and Structures

## **key\_assign\_private** Kernel Service Purpose

Requests a private kernel-key assignment.

## Syntax

```
#include <sys/types.h>
#include <sys/skeys.h>
#include <sys/kerrno.h>
```

```
kerrno_t key_assign_private (id, instance, flags, kkey)
char *id;
long instance;
unsigned long flags;
kkey_t *kkey;
```

## Parameters



Item	Description
<i>id</i>	Specifies a null-terminated string. The <code>kkey_assign_private</code> kernel service uses the string value to assign a private key. This normally contains a load module name associated with the calling kernel subsystem, but you can specify any unique string.
<i>instance</i>	Specifies a unique number for each private key requested by a subsystem. This must be an integer value starting from 0 and increases with each kernel-key requested.
<i>flags</i>	You must specify this parameter to zero.
<i>kkey</i>	Contains the returned assigned kernel key. The valid pointer must be a 4-byte aligned address ( <code>kkey_t</code> 's natural alignment).

## Description

The `kkey_assign_private` kernel service assigns a private kernel key to the caller. Private kernel keys are used to limit data accessibility by external kernel code. The `kkey_assign_private` kernel service distributes requests for private kernel keys among a predetermined range (from `KKEY_PRIVATE1` to `KKEY_PRIVATE32`). The intention is to perform a uniform distribution on behalf of requests by multiple kernel subsystems. The assignment is made based on the *id* and *instance* parameters and might return the same private key to multiple callers. It might also return the same private key when the instance number is different.

The `kkey_assign_private` kernel service does not perform a resource allocation. It only provides a recommended kernel key to use for data protection.

## Execution Environment

The `kkey_assign_private` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful completion.
<code>EINVAL_KKEY_ASSIGN_PRIVATE</code>	Indicates that the parameter or execution environment is not valid.

## kkeyset\_add\_key Kernel Service

### Purpose

Adds a kernel key to a kernel keyset.

### Syntax

```
#include <sys/kerrno.h>
#include <sys/skeys.h>
```

```
kerrno_t kkeyset_add_key (set, key, flags)
kkeyset_t set;
kkey_t key;
unsigned long flags;
```

### Parameters

Item	Description
<i>set</i>	Specifies the kernel keyset to which the <code>kkeyset_add_key</code> kernel service will add a key.
<i>key</i>	Specifies the kernel key to add.
<i>flags</i>	You can specify the <i>flags</i> parameter to one of the following values:
<b>KA_READ</b>	Specifies that the read access for the key is to be added.
<b>KA_WRITE</b>	Specifies that the write access for the key is to be added.
<b>KA_RW</b>	Specifies that both the read access and the write access are to be added. This is equivalent to the value of <code>KA_READ   KA_WRITE</code> .

## Description

The `kkeyset_add_key` kernel service adds a single kernel key specified by the *key* parameter to the kernel keyset specified by the *set* parameter. You must specify the *flags* parameter to control the read or write authority.

## Execution Environment

The `kkeyset_add_key` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful completion.
EINVAL_KKEYSET_ADD_KEY	Indicates that the parameter or execution environment is not valid.

## `kkeyset_add_set` Kernel Service Purpose

Adds members of one kernel keyset to an existing kernel keyset.

## Syntax

```
#include <sys/kerrno.h>
#include <sys/skeys.h>
```

```
kerrno_t kkeyset_add_set (set, addset)
kkeyset_t set;
kkeyset_t addset;
```

## Parameters

Item	Description
<i>set</i>	Specifies an existing kernel keyset. This set contains the resulting union on completion.
<i>addset</i>	Specifies the kernel keyset to add.

## Description

The `kkeyset_add_set` kernel service adds a kernel keyset specified by the *addset* parameter to the kernel keyset specified by the *set* parameter.

## Execution Environment

The `kkeyset_add_set` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful completion.
EINVAL_KKEYSET_ADD_SET	Indicates that the parameter or execution environment is not valid.

## kkeyset\_create Kernel Service Purpose

Creates and initializes a kernel keyset.

### Syntax

```
#include <sys/kerrno.h>
#include <sys/skeys.h>
```

```
kerrno_t kkeyset_create (set)
kkeyset_t *set;
```

### Parameters

Item	Description
<i>set</i>	Contains the returned newly-created keyset.

### Description

The **kkeyset\_create** kernel service creates a new (empty) kernel keyset. You can add or remove the access to an individual or groups of kernel keys using the **kkeyset\_add\_key**, **kkeyset\_remove\_key**, **kkeyset\_add\_set**, and **kkeyset\_remove\_set** kernel services.

**Important:** The **kkeyset\_create** kernel service allocates hidden kernel resources. You must release these resources using the **kkeyset\_delete** kernel service when the kernel keyset is no longer in use. When creating a new set, the caller of the **kkeyset\_create** kernel service must initialize the storage that will contain the returned kernel keyset (*set*) to the value of **KKEYSET\_INVALID**.

### Execution Environment

The **kkeyset\_create** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful completion.
ENOMEM_KKEYSET_CREATE	Indicates that the available memory is not sufficient to satisfy the request.
EINVAL_KKEYSET_CREATE	Indicates that the parameter or execution environment is not valid.

### Related reference:

“**kkeyset\_add\_key** Kernel Service” on page 273

“**kkeyset\_remove\_key** Kernel Service” on page 276

“**kkeyset\_delete** Kernel Service”

## kkeyset\_delete Kernel Service Purpose

Deletes a kernel keyset.

## Syntax

```
#include <sys/kernno.h>
#include <sys/skeys.h>
```

```
kernno_t kkeyset_delete (set)
kkeyset_t set;
```

## Parameters

Item	Description
<i>set</i>	Specifies the keyset to be destroyed.

## Description

The **kkeyset\_delete** kernel service destroys a kernel keyset. The kernel service releases the hidden resources associated with this keyset.

## Execution Environment

The **kkeyset\_delete** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful completion.
EINVAL_KKEYSET_DELETE	Indicates that the parameter or execution environment is not valid.

## kkeyset\_remove\_key Kernel Service Purpose

Removes a kernel key from a kernel keyset.

## Syntax

```
#include <sys/kernno.h>
#include <sys/skeys.h>
```

```
kernno_t kkeyset_remove_key (set, key, flags)
kkeyset_t set;
kkey_t key;
unsigned long flags;
```

## Parameters

Item	Description
<i>set</i>	Specifies the kernel keyset from which the <b>kkeyset_remove_key</b> kernel service will remove a key.
<i>key</i>	Specifies the kernel key to remove.
<i>flags</i>	You can specify the <i>flags</i> parameter to one of the following values:  <b>KA_READ</b> Specifies that the read access for the key is to be removed.  <b>KA_WRITE</b> Specifies that the write access for the key is to be removed.  <b>KA_RW</b> Specifies that both the read access and the write access are to be removed. This is equivalent to the value of <b>KA_READ</b>   <b>KA_WRITE</b> .

## Description

The `kkeyset_remove_key` kernel service removes a single kernel key specified by the `key` parameter from the kernel keyset specified by the `set` parameter. You must specify the `flags` parameter to control the read or write authority.

## Execution Environment

The `kkeyset_remove_key` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful completion.
<code>EINVAL_KKEYSET_REMOVE_KEY</code>	Indicates that the parameter or execution environment is not valid.

## `kkeyset_remove_set` Kernel Service

### Purpose

Removes members of one kernel keyset from an existing kernel keyset.

### Syntax

```
#include <sys/kerrno.h>
#include <sys/skeys.h>
```

```
kerrno_t kkeyset_remove_set (set, removeset)
kkeyset_t set;
kkeyset_t removeset;
```

### Parameters

Item	Description
<code>set</code>	Specifies the kernel keyset from which the <code>kkeyset_remove_set</code> kernel service will remove a keyset.
<code>removeset</code>	Specifies the kernel keyset to remove.

## Description

The `kkeyset_remove_set` kernel service removes a kernel keyset specified by the `removeset` parameter from the kernel keyset specified by the `set` parameter.

## Execution Environment

The `kkeyset_remove_set` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful completion.
EINVAL_KKEYSET_REMOVE_SET	Indicates that the parameter or execution environment is not valid.

## **kkeyset\_to\_hkeyset** Kernel Service

### **Purpose**

Computes the hardware keyset associated with a kernel keyset.

### **Syntax**

```
#include <sys/kernno.h>
#include <sys/skeys.h>
```

```
kernno_t kkeyset_to_hkeyset (kkeyset, hkeyset)
kkeyset_t kkeyset;
hkeyset_t *hkeyset;
```

### **Parameters**

Item	Description
<i>kkeyset</i>	Specifies the input kernel keyset to be mapped.
<i>hkeyset</i>	Specifies the hardware keyset that is mapped to. The valid pointer must be an 8-byte aligned address.

### **Description**

The `kkeyset_to_hkeyset` kernel service maps a kernel keyset to its associated hardware keyset.

### **Execution Environment**

The `kkeyset_to_hkeyset` kernel service can be called from the process environment only.

### **Return Values**

Item	Description
0	Indicates a successful completion.
EINVAL_KKEYSET_TO_HKEYSET	Indicates that the parameter or execution environment is not valid.

## **lpar\_get\_info** Kernel Service

### **Purpose**

Retrieves the calling partition's characteristics.

### **Syntax**

```
#include <sys/dr.h>
```

```
int lpar_get_info (command, lparinfo, bufsize)
int command;
void *lparinfo;
size_t bufsize;
```

### **Parameters**

Item	Description
<i>command</i>	Specifies whether the user wants <b>format1</b> , <b>format2</b> , or <b>processor module</b> details.
<i>lparinfo</i>	Pointer to the user-allocated buffer that is passed in.
<i>bufsize</i>	Size of the buffer that is passed in.

## Description

The **klpar\_get\_info** kernel service retrieves LPAR and Micro-Partitioning attributes of both low-frequency use and high-frequency use and also retrieves processor module information. Because the low-frequency attributes, as defined in the **lpar\_info\_format1\_t** structure, are static in nature, a reboot is required to effect any change. The high-frequency attributes, as defined in the **lpar\_info\_format2\_t** structure, can be changed dynamically while the partition is running. The signature of this kernel service, its parameter types, and the order of the member fields in both the **lpar\_info\_format1\_t** and **lpar\_info\_format2\_t** structures are specific to the AIX platform. If you requests **processor module** information, the kernel service provides this information as an array of **proc\_module\_info\_t** structures. To obtain this information, the caller must provide a buffer of the exact length to accommodate one **proc\_module\_info\_t** structure for each module type. You can obtain the module count using the **NUM\_PROC\_MODULE\_TYPES** command. The module count is in the form of a **uint64\_t** type. Processor module information is reported for the entire system. This information is available on POWER6® and later systems.

To see the complete structures of **lpar\_info\_format1\_t**, **lpar\_info\_format2\_t**, and **proc\_module\_info\_t**, refer to the **dr.h** header file.

## Return Values

Upon success, the **klpar\_get\_info** kernel service returns a value of 0. Upon failure, the **klpar\_get\_info** kernel service returns an error code.

## Error Codes

Item	Description
EINVAL	Invalid input parameter.
ENOSYS	The hardware or the firmware level does not support this operation.
ENOTSUP	The platform does not support this operation.

## Related information:

**lpar\_get\_info** subroutine

## kmod\_entrypt Kernel Service Purpose

Returns a function pointer to a kernel module's entry point.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/ldr.h>
```

```
void (*(kmod_entrypt ( kmid, flags)))( )
mid_t kmid;
uint flags;
```

## Parameters

Item	Description		
<i>kmid</i>	Specifies the kernel module ID of the object file for which the entry point is requested. This parameter is the kernel module ID returned by the <b>kmod_load</b> kernel service.		
<i>flags</i>	Flag specifying entry point options. The following flag is defined: <table border="0"> <tr> <td style="padding-left: 20px;">0</td> <td>Returns a function pointer to the specified module's entry point as specified in the module header.</td> </tr> </table>	0	Returns a function pointer to the specified module's entry point as specified in the module header.
0	Returns a function pointer to the specified module's entry point as specified in the module header.		

## Description

The **kmod\_entrypt** kernel service obtains a function pointer to a specified module's entry point. This function pointer is typically used to invoke a routine in the module for initializing or terminating its functions. Initialization and termination occurs after loading and before unloading. The module for which the entry point is requested is specified by the kernel module ID represented by the *kmid* parameter.

## Execution Environment

The **kmod\_entrypt** kernel service can be called from the process environment only.

## Return Values

A nonnull function pointer indicates a successful completion. This function pointer contains the module's entry point. A null function pointer indicates an error.

### Related reference:

“kmod\_load Kernel Service”

### Related information:

Kernel Extension and Device Driver Management Kernel Services

## kmod\_load Kernel Service

### Purpose

Loads an object file into the kernel or queries for an object file already loaded.

### Syntax

```
#include <sys/ldr.h>
#include <sys/types.h>
#include <sys/errno.h>
```

```
int kmod_load (pathp,
flags, libpathp, kmidp)
caddr_t pathp;
uint flags;
caddr_t
libpathp;
mid_t * kmidp;
```

### Parameters



<b>Item</b>	<b>Description</b>
<i>pathp</i>	Points to a character string containing the path-name of the object file to load or query.
<i>flags</i>	Specifies a set of loader flags describing which loader options to invoke. The following flags are defined: <p><b>LD_USRPATH</b> The character strings pointed to by the <i>pathp</i> and <i>libpathp</i> parameters are in user address space. If the <b>LD_USRPATH</b> flag is not set, the character strings are assumed to be in kernel, or system, space.</p> <p><b>LD_KERNELEX</b> Puts this object file's exported symbols into the <b>/usr/lib/boot/unix</b> name space. Additional object files loaded due to symbol resolution for the specified file do not have their exported symbols placed in kernel name space.</p> <p><b>LD_SINGLELOAD</b> When this flag is set, the object file specified by the <i>pathp</i> parameter is loaded into the kernel only if an object file with the same path-name has not already been loaded. If an object file with the same path-name has already been loaded, its module ID is returned (using the <i>kmidp</i> parameter) and its load count incremented. If the object file is not yet loaded, this service performs the load as if the flag were not set.</p> <p>This option is useful in supporting global kernel routines where only one copy of the routine and its data can be present. Typically, routines that export symbols to be added to kernel name space are of this type.</p> <p><b>Note:</b> A path-name comparison is done to determine whether the same object file has already been loaded. This service will erroneously load a new copy of the object file into the kernel if the path-name to the object file is expressed differently than it was on a previous load request.</p> <p>If neither this flag nor the <b>LD_QUERY</b> flag is set, this service loads a new copy of the object file into the kernel. This occurs even if other copies of the object file have previously been loaded.</p> <p><b>LD_QUERY</b> This flag specifies that a query operation will determine if the object file specified by the <i>pathp</i> parameter is loaded. If not loaded, a kernel module ID of 0 is returned using the <i>kmidp</i> parameter. Otherwise, the kernel module ID assigned to the object file is returned.</p> <p>If multiple instances of this file have been loaded into the kernel, the kernel module ID of the most recently loaded object file is returned.</p> <p>The <i>libpathp</i> parameter is not used for this option.</p> <p><b>Note:</b> A path-name comparison is done to determine whether the same object file has been loaded. This service will erroneously return a not loaded condition if the path-name to the object file is expressed differently than it was on a previous load request.</p> <p>If this flag is set, no object file is loaded and the <b>LD_SINGLELOAD</b> and <b>LD_KERNELEX</b> flags are ignored, if set.</p>
<i>libpathp</i>	Points to a character string containing the search path to use for finding object files required to complete symbol resolution for this load. If the parameter is null, the search path is set from the specification in the object file header for the object file specified by the <i>pathp</i> parameter.
<i>kmidp</i>	Points to an area where the kernel module ID associated with this load of the specified module is to be returned. The data in this area is not valid if the <b>kmod_load</b> service returns a nonzero return code.

## Description

The **kmod\_load** kernel service loads into the kernel a kernel extension object file specified by the *pathp* parameter. This service returns a kernel module ID for that instance of the module.

You can specify flags to request a single load, which ensures that only one copy of the object file is loaded into the kernel. An additional option is simply to query for a given object file (specified by path-name). This allows the user to determine if a module is already loaded and then access its assigned kernel module ID.

The **kmod\_load** service also provides load-time symbol resolution of the loaded module's imported symbols. The **kmod\_load** service loads additional kernel object modules if required for symbol resolution.

## Loader Symbol Binding Support

Symbols imported from the kernel name space are resolved with symbols that exist in the kernel name space at the time of the load. (Symbols are imported from the kernel name space by specifying the `#!/unix` character string as the first field in an import list at link-edit time.)

Kernel modules can also import symbols from other kernel object modules. These other kernel object modules are loaded along with the specified object module if they are needed to resolve the imported symbols.

Any symbols exported by the specified kernel object module are added to the kernel name space if the *flags* parameter has the `LD_KERNELEX` flag set. This makes the symbols available to other subsequently loaded kernel object modules. Kernel object modules loaded on behalf of the specified kernel object module (to resolve imported symbols) do not have their exported symbols added to the kernel name space.

Kernel export symbols specified (at link-edit time) with the `SYSCALL` keyword in the primary module's export list are added to the system call table. These kernel export symbols are available to application programs as system calls.

### **Finding Shared Object Modules for Resolving Symbol References**

The search path search string is taken from the module header of the object module specified by the *pathp* parameter if the *libpathp* parameter is null. The module header of the object module specified by the *pathp* parameter is used.

If the module header contains an unqualified base file name for the symbol (no / [slash] characters in the name), a search string is used to find the location of the shared object module required to resolve the import. This search string can be taken from one of two places. If the *libpathp* parameter on the call to the `kmod_load` service is not null, then it points to a character string specifying the search path to be used. However, if the *libpathp* parameter is null, then the search path is to be taken from the module header for the object module specified by the *pathp* parameter.

The search path specification found in object modules loaded to resolve imported symbols is not used. The kernel loader service does not support deferred symbol resolution. The load of the kernel module is terminated with an error if any imported symbols cannot be resolved.

### **Execution Environment**

The `kmod_load` kernel service can be called from the process environment only.

### **Return Values**

If the object file is loaded without error, the module ID is returned in the location pointed to by the *kmidp* parameter and the return code is set to 0.

### **Error Codes**

If an error results, the module is not loaded, and no kernel module ID is returned. The return code is set to one of the following return values:

Return Value	Description
EACCES	Indicates that an object module to be loaded is not an ordinary file or that the mode of the object module file denies read-only access.
EACCES	Search permission is denied on a component of the path prefix.
EFAULT	Indicates that the calling process does not have sufficient authority to access the data area described by the <i>pathp</i> or <i>libpathp</i> parameters when the <b>LD_USRPATH</b> flag is set. This error code is also returned if an I/O error occurs when accessing data in this area.
ENOEXEC	Indicates that the program file has the appropriate access permission, but has an XCOFF indicator that is not valid in its header. The <b>kmod_load</b> kernel service supports loading of XCOFF (Extended Common Object File Format) object files only. This error code is also returned if the loader is unable to resolve an imported symbol.
EINVAL	Indicates that the program file has a valid XCOFF indicator in its header, but the header is either damaged or incorrect for the machine on which the file is to be loaded.
ENOMEM	Indicates that the load requires more kernel memory than allowed by the system-imposed maximum.
ETXTBSY	Indicates that the object file is currently open for writing by some process.
ENOTDIR	Indicates that a component of the path prefix is not a directory.
ENOENT	Indicates that no such file or directory exists or the path-name is null.
ESTALE	Indicates that the caller's root or current directory is located in a virtual file system that has been unmounted.
ELOOP	Indicates that too many symbolic links were encountered in translating the <i>path</i> or <i>libpathp</i> parameter.
ENAMETOOLONG	Indicates that a component of a path-name exceeded 255 characters, or an entire path-name exceeded 1023 characters.
EIO	Indicates that an I/O error occurred during the operation.

**Related reference:**

“kmod\_unload Kernel Service”

**Related information:**

kmod\_util subroutine

Kernel Extension and Device Driver Management Kernel Services

## kmod\_unload Kernel Service

### Purpose

Unloads a kernel object file.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/ldr.h>
```

```
int kmod_unload ( kmid
, flags)
mid_t kmid;
uint flags;
```

### Parameters

Item	Description
<i>kmid</i>	Specifies the kernel module ID of the object file to be unloaded. This kernel module ID is returned when using the <b>kmod_load</b> kernel service.
<i>flags</i>	Flags specifying unload options. The following flag is defined: <ul style="list-style-type: none"> <li>0 Unloads the object module specified by its <i>kmid</i> parameter and any object modules that were loaded as a result of loading the specified object file if this file is not still in use.</li> </ul>

## Description

The **kmod\_unload** kernel service unloads a previously loaded kernel extension object file. The object to be unloaded is specified by the *kmid* parameter. Upon successful completion, the following objects are unloaded or marked *unload pending*:

- The specified object file
- Any imported kernel object modules that were loaded as a result of the loading of the specified module

Users of these exports or system calls are modules bound to this module's exported symbols. If there are no users of any of the module's kernel exports or system calls, the module is immediately unloaded. If there are users of this module, the module is not unloaded but marked *unload pending*.

Marking a module *unload pending* removes the module's exported symbols from the kernel name space. Any system calls exported by this module are also removed. This prohibits new users of these symbols. The module is unloaded only when all current users have been unloaded.

If the unload is successfully completed or marked *pending*, a value of 0 is returned. When an error occurs, the specified module and any imported modules are not unloaded. A nonzero return value indicates the error.

## Execution Environment

The **kmod\_unload** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates successful completion.
EINVAL	Indicates that the <i>kmid</i> parameter, which specifies the kernel module, is not valid or does not correspond to a currently loaded module.
EBUSY	The <i>kmid</i> parameter specifies a kernel extension that is still intercepting system calls.

### Related reference:

“kmod\_load Kernel Service” on page 280

### Related information:

kmod\_util subroutine

Kernel Extension and Device Driver Management Kernel Services

## kmod\_util Kernel Service

### Purpose

Registers routines to be called before and after specified system calls are invoked.

### Syntax

```
#include <sys/sysconfig.h>
```

```
int kmod_util ( flags, buffer, blen)
int flags;
void * buffer;
long blen;
```

## Parameters

Item	Description
<i>flags</i>	Specifies the operation. Valid values are <b>KU_INTERCEPT</b> , <b>KU_INTERCEPT_STOP</b> , and <b>KU_INTERCEPT_CANCEL</b> .
<i>buffer</i>	Points to a buffer containing a system call interception header and an array of system call interception structures.
<i>blen</i>	Specifies the length of the buffer.

## Description

The **kmod\_util** kernel service allows system calls to be intercepted. Routines called **pre-sc** functions are specified to be called before the intercepted system call. Routines called **post-sc** functions are specified to be called after the intercepted system call. In addition, a **pre-sc** function is allowed to abort the system call, providing its own return value and preventing subsequent **pre-sc** functions and the system call itself from being called. Similarly, each **post-sc** function may examine and alter the return value. If a system call does not return (e.g., **thread\_terminate**), **post-sc** functions are not called.

For each intercepted system call, either a **pre-sc** function, a **post-sc** function, or both, must be specified. If a **pre-sc** function and **post-sc** function are registered for the same system call in the same **kmod\_util** invocation, they are considered paired. All **pre-sc** and **post-sc** functions specified in a **kmod\_util** call must be defined in the same kernel extension as the caller of the **kmod\_util** kernel service. Other kernel extensions, however, can intercept the same system calls. The most recently registered **pre-sc** function is called first, and its paired **post-sc** function is called last.

The interception of a system call is implemented so that all calls to the system call are intercepted, even for existing processes.

It may be necessary to prevent the interception of certain system calls to avoid destabilizing the system. A future version or release of the **kmod\_util** kernel service may prevent the interception of additional system calls, and such a change will not be considered a violation of binary compatibility.

The prototype of a **pre-sc** function is

```
int pre_sc(uintptr_t *rc, void *parms, uintptr_t cookie, void *buffer);
```

where *parms* is a pointer to the parameters of the system call, *cookie* is an opaque value specified by the caller of **kmod\_util**, *buffer* is a scratch 128-byte buffer for use by the **pre-sc** function and its paired **post-sc** function.

If the **pre-sc** function returns non-zero, the system call is aborted. The *rc* parameter is the address where an alternate return value can be specified. Subsequent **pre\_sc** functions are not called, nor is the system call. For **pre-sc** functions already called, their paired **post-sc** functions are called.

The prototype of a **post-sc** function is

```
void post_sc(uintptr_t *rc, void *parms, uintptr_t cookie,
            void *buffer);
```

The parameters of the **post-sc** function are the same as those of the **pre-sc** function. In particular, the *buffer* parameter is the same buffer that was passed to the paired **pre-sc** function. The return value can be modified by a **post-sc** function.

For calls to the **kmod\_util** kernel service, the buffer contains a header and an array of elements about system calls to be intercepted. The layout of these structures is defined in `<sys/sysconfig.h>`.

An array element is ignored if the **KU\_IGNORE** flag is set in the **kue\_iflag** field. Otherwise, each array element in the input buffer is validated, and if any errors are found, the entire call fails without any partial execution.

### Intercepting System Calls

Calls to **kmod\_util()** with the **KU\_INTERCEPT** flag initiate system call interception.

### Stopping System Call Interception

Calls to **kmod\_util()** with the **KU\_INTERCEPT\_STOP** flag suspend the interception of the specified system calls. If a **pre-sc** function has already been called for a specified system call, its paired **post-sc** function will still be called, but future calls to the system call will not invoke either the **pre-sc** or **post-sc** function. It is not valid to stop interception of a system call that was not originally intercepted by the calling kernel extension.

If the interception of a system call has been suspended, it may be resumed by calling the **kmod\_util()** function with the **KU\_INTERCEPT** flag, as long as the same values are specified, such as the **pre-sc** and **post-sc** functions.

### Cancelling System Call Interception

System call interception can be cancelled by specifying the **KU\_INTERCEPT\_CANCEL** flag. When interception is cancelled, the **post-sc** function is not called even if its paired **pre-sc** function was called. It is not valid to cancel interception of a system call that was not originally intercepted by the calling kernel extension, but interception can be cancelled without first stopping interception.

Once interception of a system call has been cancelled, it can be intercepted anew by calling the **kmod\_util()** function with the **KU\_INTERCEPT** flag. Different **pre-sc** and **post-sc** functions can be specified in this case.

## Return Values

If the specified operations can be enacted for all specified system calls, 0 is returned. Otherwise, a non-zero value is returned and no change in the state of system call interception occurs. If an error occurs because of a validation error in a particular array element, the **kue\_oflags** field usually identifies the error in more detail.

## Error Codes

If an error results, one of the following error values is returned:

Return Value	Description
EINVAL	The flags parameter is not <b>KU_INTERCEPT</b> , <b>KU_INTERCEPT_STOP</b> , nor <b>KU_INTERCEPT_CANCEL</b> .  The fields in the header are invalid or the <i>blen</i> parameter is not consistent with the number of array elements.  The buffer was invalid. For <b>KU_INTERCEPT</b> , at least one of the <b>pre-sc</b> and <b>post-sc</b> must be supplied for each system call to be intercepted. All <b>pre-sc</b> and <b>post-sc</b> functions must be in the same kernel extension as the caller of <b>kmod_util()</b> .
EBUSY	A request was made to intercept a system call that was already being intercepted.
ENOENT	A request was made to stop or cancel interception of a system call that was not being intercepted.
ENOMEM	Memory could not be allocated to satisfy the request.
ENOTSUPP	One of the specified system calls is not allowed to be intercepted.

### Related reference:

“kmod\_load Kernel Service” on page 280

“kmod\_unload Kernel Service” on page 283

### Related information:

Kernel Extension and Device Driver Management Kernel Services

## kmsgctl Kernel Service

### Purpose

Provides message-queue control operations.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/ipc.h>
#include <sys/msg.h>
```

```
int kmsgctl ( msqid, cmd, buf)
int msqid, cmd;
struct msqid_ds *buf;
```

### Parameters

Item	Description
<i>msqid</i>	Specifies the message queue ID, which indicates the message queue for which the control operation is being requested for.
<i>cmd</i>	Specifies which control operation is being requested. There are three valid commands.
<i>buf</i>	Points to the <b>msqid_ds</b> structure provided by the caller of the <b>kmsgctl</b> service. Data is obtained either from this structure or from status returned in this structure, depending on the <i>cmd</i> parameter. The <b>msqid_ds</b> structure is defined in the <code>/usr/include/sys/msg.h</code> file.

### Description

The **kmsgctl** kernel service provides a variety of message-queue control operations as specified by the *cmd* parameter. The **kmsgctl** kernel service provides the same functions for user-mode processes in kernel mode as the **msgctl** subroutine performs for kernel processes or user-mode processes in user mode. The **kmsgctl** service can be called by a user-mode process in kernel mode or by a kernel process. A kernel process can also call the **msgctl** subroutine to provide the same function.

The following three commands can be specified with the *cmd* parameter:

Item	Description
IPC_STAT	Sets only documented fields. See the <b>msgctl</b> subroutine.
IPC_SET	Sets the value of the following fields of the data structure associated with the <i>msqid</i> parameter to the corresponding values found in the structure pointed to by the <i>buf</i> parameter: <ul style="list-style-type: none"><li>msg_perm.uid</li><li>msg_perm.gid</li><li>msg_perm.mode (only the low-order 9 bits)</li><li>msg_qbytes</li></ul>
IPC_RMID	To perform the <b>IPC_SET</b> operation, the current process must have an effective user ID equal to the value of the msg_perm.uid or msg_perm.cuid field in the data structure associated with the <i>msqid</i> parameter. To raise the value of the msg_qbytes field, the calling process must have the appropriate system privilege. Removes from the system the message-queue identifier specified by the <i>msqid</i> parameter. This operation also destroys both the message queue and the data structure associated with it. To perform this operation, the current process must have an effective user ID equal to the value of the msg_perm.uid or msg_perm.cuid field in the data structure associated with the <i>msqid</i> parameter.

## Execution Environment

The `kmsgctl` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates successful completion.
EINVAL	Indicates either <ul style="list-style-type: none"><li>• The identifier specified by the <code>msgid</code> parameter is not a valid message queue identifier.</li><li>• The command specified by the <code>cmd</code> parameter is not a valid command.</li></ul>
EACCES	The command specified by the <code>cmd</code> parameter is equal to <code>IPC_STAT</code> and read permission is denied to the calling process.
EPERM	The command specified by the <code>cmd</code> parameter is equal to <code>IPC_RMID</code> , <code>IPC_SET</code> , and the effective user ID of the calling process is not equal to that of the value of the <code>msg_perm.uid</code> field in the data structure associated with the <code>msgid</code> parameter.
EPERM	Indicates the following conditions: <ul style="list-style-type: none"><li>• The command specified by the <code>cmd</code> parameter is equal to <code>IPC_SET</code>.</li><li>• An attempt is being made to increase to the value of the <code>msg_qbytes</code> field, but the calling process does not have the appropriate system privilege.</li></ul>

### Related information:

`msgctl` subroutine

Message Queue Kernel Services

Understanding System Call Execution

## kmsgget Kernel Service

### Purpose

Obtains a message queue identifier.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/stat.h>
#include <sys/ipc.h>
#include <sys/msg.h>
```

```
int kmsgget ( key, msgflg, msgid)
```

```
key_t key;
```

```
int msgflg;
```

```
int *msgid;
```

### Parameters

Item	Description
<i>key</i>	Specifies either a value of <code>IPC_PRIVATE</code> or an IPC key constructed by the <code>ftok</code> subroutine (or a similar algorithm).



Item	Description
<i>msgflg</i>	Specifies that the <i>msgflg</i> parameter is constructed by logically ORing one or more of these values: <ul style="list-style-type: none"> <li><b>IPC_CREAT</b> Creates the data structure if it does not already exist.</li> <li><b>IPC_EXCL</b> Causes the <b>kmsgget</b> kernel service to fail if <b>IPC_CREAT</b> is also set and the data structure already exists.</li> <li><b>S_IRUSR</b> Permits the process that owns the data structure to read it.</li> <li><b>S_IWUSR</b> Permits the process that owns the data structure to modify it.</li> <li><b>S_IRGRP</b> Permits the process group associated with the data structure to read it.</li> <li><b>S_IWGRP</b> Permits the process group associated with the data structure to modify it.</li> <li><b>S_IROTH</b> Permits others to read the data structure.</li> <li><b>S_IWOTH</b> Permits others to modify the data structure.</li> </ul> <p>The values that begin with <b>S_I...</b> are defined in the <code>/usr/include/sys/stat.h</code> file. They are a subset of the access permissions that apply to files.</p>
<i>msqid</i>	A reference parameter where a valid message-queue ID is returned if the <b>kmsgget</b> kernel service is successful.

## Description

The **kmsgget** kernel service returns the message-queue identifier specified by the *msqid* parameter associated with the specified *key* parameter value. The **kmsgget** kernel service provides the same functions for user-mode processes in kernel mode as the **msgget** subroutine performs for kernel processes or user-mode processes in user mode. The **kmsgget** service can be called by a user-mode process in kernel mode or by a kernel process. A kernel process can also call the **msgget** subroutine to provide the same function.

## Execution Environment

The **kmsgget** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates successful completion. The <i>msqid</i> parameter is set to a valid message-queue identifier.

If the **kmsgget** kernel service fails, the *msqid* parameter is not valid and the return code is one of these four values:

Item	Description
EACCES	Indicates that a message queue ID exists for the <i>key</i> parameter but operation permission as specified by the <i>msgflg</i> parameter cannot be granted.
ENOENT	Indicates that a message queue ID does not exist for the <i>key</i> parameter and the <b>IPC_CREAT</b> command is not set.
ENOSPC	Indicates that a message queue ID is to be created but the system-imposed limit on the maximum number of allowed message queue IDs systemwide will be exceeded.
EEXIST	Indicates that a message queue ID exists for the value specified by the <i>key</i> parameter, and both the <b>IPC_CREAT</b> and <b>IPC_EXCL</b> commands are set.

## Related information:

msgget subroutine

## kmsgrcv Kernel Service Purpose

Reads a message from a message queue.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/ipc.h>
#include <sys/msg.h>
```

```
int kmsgrcv
(msqid, msgp, msgsz,
msgtyp, msgflg, flags, bytes)
int msqid;
struct msgxbuf * msgp;
    or struct msgbuf *msgp;
int msgsz;
mtyp_t msgtyp;
int msgflg;
int flags;
ssize_t * bytes;
```

### Parameters

Item	Description
<i>msqid</i>	Specifies the message queue from which to read.
<i>msgp</i>	Points to either an <b>msgxbuf</b> or an <b>msgbuf</b> structure where the message text is placed. The type of structure pointed to is determined by the values of the <i>flags</i> parameter. These structures are defined in the <code>/usr/include/sys/msg.h</code> file.
<i>msgsz</i>	Specifies the maximum number of bytes of text to be received from the message queue. The received message is truncated to the size specified by the <i>msgsz</i> parameter if the message is longer than this size and <b>MSG_NOERROR</b> is set in the <i>msgflg</i> parameter. The truncated part of the message is lost and no indication of the truncation is given to the calling process.
<i>msgtyp</i>	Specifies the type of message requested as follows: <ul style="list-style-type: none"><li>• If the <i>msgtyp</i> parameter is equal to 0, the first message on the queue is received.</li><li>• If the <i>msgtyp</i> parameter is greater than 0, the first message of the type specified by the <i>msgtyp</i> parameter is received.</li><li>• If the <i>msgtyp</i> parameter is less than 0, the first message of the lowest type that is less than or equal to the absolute value of the <i>msgtyp</i> parameter is received.</li></ul>
<i>msgflg</i>	Specifies a value of 0, or is constructed by logically ORing one of several values: <b>MSG_NOERROR</b> Truncates the message if it is longer than the number of bytes specified by the <i>msgsz</i> parameter. <b>IPC_NOWAIT</b> Specifies the action to take if a message of the desired type is not on the queue: <ul style="list-style-type: none"><li>• If <b>IPC_NOWAIT</b> is set, then the <b>kmsgrcv</b> service returns an <b>ENOMSG</b> value.</li><li>• If <b>IPC_NOWAIT</b> is not set, then the calling process suspends execution until one of the following occurs:<ul style="list-style-type: none"><li>– A message of the desired type is placed on the queue.</li><li>– The message queue ID specified by the <i>msqid</i> parameter is removed from the system. When this occurs, the <b>kmsgrcv</b> service returns an <b>EIDRM</b> value.</li><li>– The calling process receives a signal that is to be caught. In this case, a message is not received and the <b>kmsgrcv</b> service returns an <b>EINTR</b> value.</li></ul></li></ul>

Item	Description
<i>flags</i>	Specifies a value of 0 if a normal message receive is to be performed. If an extended message receive is to be performed, this flag should be set to an <b>XMSG</b> value. With this flag set, the <b>kmsgrcv</b> service functions as the <b>msgxrcv</b> subroutine would. Otherwise, the <b>kmsgrcv</b> service functions as the <b>msgrcv</b> subroutine would.
<i>bytes</i>	Specifies a reference parameter. This parameter contains the number of message-text bytes read from the message queue upon return from the <b>kmsgrcv</b> service.

If the message is longer than the number of bytes specified by the *msgsz* parameter bytes but **MSG\_NOERROR** is not set, then the **kmsgrcv** kernel service fails and returns an **E2BIG** return value.

## Description

The **kmsgrcv** kernel service reads a message from the queue specified by the *msqid* parameter and stores the message into the structure pointed to by the *msgp* parameter. The **kmsgrcv** kernel service provides the same functions for user-mode processes in kernel mode as the **msgrcv** and **msgxrcv** subroutines perform for kernel processes or user-mode processes in user mode.

The **kmsgrcv** service can be called by a user-mode process in kernel mode or by a kernel process. A kernel process can also call the **msgrcv** and **msgxrcv** subroutines to provide the same functions.

## Execution Environment

The **kmsgrcv** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
<b>EINVAL</b>	Indicates that the ID specified by the <i>msqid</i> parameter is not a valid message queue ID.
<b>EACCES</b>	Indicates that operation permission is denied to the calling process.
<b>EINVAL</b>	Indicates that the value of the <i>msgsz</i> parameter is less than 0.
<b>E2BIG</b>	Indicates that the message text is greater than the maximum length specified by the <i>msgsz</i> parameter and <b>MSG_NOERROR</b> is not set.
<b>ENOMSG</b>	Indicates that the queue does not contain a message of the desired type and <b>IPC_NOWAIT</b> is set.
<b>EINTR</b>	Indicates that the <b>kmsgrcv</b> service received a signal.
<b>EIDRM</b>	Indicates that the message queue ID specified by the <i>msqid</i> parameter has been removed from the system.

### Related information:

msgrcv subroutine

msgxrcv subroutine

Message Queue Kernel Services

## kmsgsnd Kernel Service

### Purpose

Sends a message using a previously defined message queue.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/ipc.h>
#include <sys/msg.h>
```

```
int kmsgsnd (msqid, msgp, msgsz, msgflg)
int msqid;
struct msgbuf * msgp;
int msgsz, msgflg;
```

## Parameters

Item	Description
<i>msqid</i>	Specifies the message queue ID that indicates which message queue the message is to be sent on.
<i>msgp</i>	Points to an <b>msgbuf</b> structure containing the message. The <b>msgbuf</b> structure is defined in the <code>/usr/include/sys/msg.h</code> file.
<i>msgsz</i>	Specifies the size of the message to be sent in bytes. The <i>msgsz</i> parameter can range from 0 to a system-imposed maximum.
<i>msgflg</i>	Specifies the action to be taken if the message cannot be sent for one of several reasons.

## Description

The **kmsgsnd** kernel service sends a message to the queue specified by the *msqid* parameter. The **kmsgsnd** kernel service provides the same functions for user-mode processes in kernel mode as the **msgsnd** subroutine performs for kernel processes or user-mode processes in user mode. The **kmsgsnd** service can be called by a user-mode process in kernel mode or by a kernel process. A kernel process can also call the **msgsnd** subroutine to provide the same function.

There are two reasons why the **kmsgsnd** kernel service cannot send the message:

- The number of bytes already on the queue is equal to the **msg\_qbytes** member.
- The total number of messages on all queues systemwide is equal to a system-imposed limit.

There are several actions to take when the **kmsgsnd** kernel service cannot send the message:

- If the *msgflg* parameter is set to **IPC\_NOWAIT**, then the message is not sent, and the **kmsgsnd** service fails and returns an **EAGAIN** value.
- If the *msgflg* parameter is 0, then the calling process suspends execution until one of the following occurs:
  - The condition responsible for the suspension no longer exists, in which case the message is sent.
  - The message queue ID specified by the *msqid* parameter is removed from the system. When this occurs, the **kmsgsnd** service fails and an **EIDRM** value is returned.
  - The calling process receives a signal that is to be caught. In this case, the message is not sent and the calling process resumes execution as described in the **sigaction** kernel service.

## Execution Environment

The **kmsgsnd** kernel service can be called from the process environment only.

The calling process must have write permission to perform the **kmsgsnd** operation.

## Return Values

Item	Description
0	Indicates a successful operation.
<b>EINVAL</b>	Indicates that the <i>msqid</i> parameter is not a valid message queue ID.
<b>EACCES</b>	Indicates that operation permission is denied to the calling process.
<b>EAGAIN</b>	Indicates that the message cannot be sent for one of the reasons stated previously, and the <i>msgflg</i> parameter is set to <b>IPC_NOWAIT</b> .
<b>EINVAL</b>	Indicates that the <i>msgsz</i> parameter is less than 0 or greater than the system-imposed limit.
<b>EINTR</b>	Indicates that the <b>kmsgsnd</b> service received a signal.
<b>EIDRM</b>	Indicates that the message queue ID specified by the <i>msqid</i> parameter has been removed from the system.
<b>ENOMEM</b>	Indicates that the system does not have enough memory to send the message.

### Related information:

msgsnd subroutine

Message Queue Kernel Services

## **kra\_attachrset Subroutine**

### **Purpose**

Attaches a work component to a resource set.

### **Syntax**

```
#include <sys/rset.h>
int kra_attachrset (rstype, rsid, rset, flags)
rstype_t rstype;
rsid_t rsid;
rsethandle_t rset;
unsigned int flags;
```

### **Description**

The `kra_attachrset` subroutine attaches a work component specified by the `rstype` and `rsid` parameters to a resource set specified by the `rset` parameter.

The work component is an existing process identified by the process ID or an existing kernel thread identified by the kernel thread ID (tid). A process ID or thread ID value of `RS_MYSELF` indicates the attachment applies to the current process or the current kernel thread, respectively.

The following conditions must be met to successfully attach a process to a resource set:

- The resource set must contain processors that are available in the system.
- The calling process must either have root authority or have `CAP_NUMA_ATTACH` capability.
- The calling process must either have root authority or the same effective userid as the target process.
- The target process must not contain any threads that have bindprocessor bindings to a processor.
- The resource set must be contained in (be a subset of ) the target process' partition resource set.
- The resource set must be a superset of all the thread's `rset` in the target process.

The following conditions must be met to successfully attach a kernel thread to a resource set:

- The resource set must contain processors that are available in the system.
- The calling process must either have root authority or have `CAP_NUMA_ATTACH` capability.
- The calling process must either have root authority or the same effective userid as the target process.
- The target thread must not have bindprocessor bindings to a processor.
- The resource set must be contained in (be a subset of ) the target thread's process effective and partition resource set.

If any of these conditions are not met, the attachment will fail.

Once a process is attached to a resource set, the threads in the process will only run on processors contained in the resource set. Once a kernel thread is attached to a resource set, that thread will only run on processors contained in the resource set.

The `flags` parameter can be set to indicate the policy for using the resources contained in the resource set specified in the `rset` parameter. The only supported scheduling policy is `R_ATTACH_STRSET`, which is useful only when the processors of the system are running in simultaneous multithreading mode. Processors like the POWER5 support simultaneous multithreading, where each physical processor has two execution engines, called *hardware threads*. Each hardware thread is essentially equivalent to a single CPU, and each is identified as a separate CPU in a resource set. The `R_ATTACH_STRSET` flag indicates that the process is to be scheduled with a single-threaded policy; namely, that it should be scheduled on only one hardware thread per physical processor. If this flag is specified, then all of the available

processors indicated in the resource set must be of exclusive use. A new resource set, called an *ST resource set*, is constructed from the specified resource set and attached to the process according to the following rules:

- All offline processors are ignored.
- If all the hardware threads (CPUs) of a physical processor (when running in simultaneous multithreading mode, there will be more than one active hardware thread per physical processor) are not included in the specified resource set, the other CPUs of the processor are ignored when constructing the ST resource set.
- Only one CPU (hardware thread) resource per physical processor is included in the ST resource set.

## Parameters

Item	Description
<i>rstype</i>	Specifies the type of work component to be attached to the resource set specified by the <i>rset</i> parameter. The <i>rstype</i> parameter must be the following value, defined in <b>rset.h</b> : <ul style="list-style-type: none"> <li>• R_PROCESS: existing process</li> <li>• R_THREAD: existing kernel thread</li> </ul>
<i>rsid</i>	Identifies the work component to be attached to the resource set specified by the <i>rset</i> parameter. The <i>rsid</i> parameter must be the following: <ul style="list-style-type: none"> <li>• Process ID (for <i>rstype</i> of R_PROCESS): set the <i>rsid_t at_pid</i> field to the desired process' process ID.</li> <li>• Kernel thread ID (for <i>rstype</i> of R_THREAD): set the <i>rsid_t at_tid</i> field to the desired kernel thread's thread ID.</li> </ul>
<i>rset</i>	Specifies which work component (specified by the <i>rstype</i> and <i>rsid</i> parameters) to attach to the resource set.
<i>flags</i>	Specifies the scheduling policy for the work component being attached. <p>The only supported value is R_ATTACH_STRSET value, which is only applicable if the <i>rstype</i> parameter is set to R_PROCESS. The R_ATTACH_STRSET value indicates that the process is to be scheduled with a single-threaded policy (only on one hardware thread per physical processor).</p>

## Return Values

Upon successful completion, the **kra\_attachrset** subroutine returns a 0. If unsuccessful, one or more of the following are true:

Item	Description
EINVAL	One of the following is true: <ul style="list-style-type: none"> <li>• The <i>flags</i> parameter contains an invalid value.</li> <li>• The <i>rstype</i> parameter contains an invalid type qualifier.</li> <li>• The R_ATTACH_STRSET <i>flags</i> parameter is specified and one or more processors in the <i>rset</i> parameter are not assigned for exclusive use.</li> </ul>
ENODEV	The resource set specified by the <i>rset</i> parameter does not contain any available processors, or the R_ATTACH_STRSET <i>flags</i> parameter is specified and the constructed ST resource set does not have any available processors.
ESRCH	The process or kernel thread identified by the <i>rstype</i> and <i>rsid</i> parameters does not exist.
EPERM	One of the following is true: <ul style="list-style-type: none"> <li>• If the <i>rstype</i> is R_PROCESS, either the resource set specified by the <i>rset</i> parameter is not included in the partition resource set of the process identified by the <i>rstype</i> and <i>rsid</i> parameters, or any of the thread's R_THREAD <i>rset</i> in this process is not a subset of the resource set specified by the <i>rset</i> parameter.</li> <li>• If the <i>rstype</i> is R_THREAD, the resource set specified by the <i>rset</i> parameter is not included in the target thread's process effective or partition (real) resource set.</li> <li>• The calling process has neither root authority nor CAP_NUMA_ATTACH attachment privilege.</li> <li>• The calling process has neither root authority nor the same effective user ID as the process identified by the <i>rstype</i> and <i>rsid</i> parameters.</li> <li>• The process or thread identified by the <i>rstype</i> and <i>rsid</i> parameters has one or more threads with a bindprocessor processor binding.</li> </ul>

**Related reference:**

“kra\_getrset Subroutine” on page 297

“kra\_detachset Subroutine” on page 296

#### Related information:

Exclusive use processor resource sets

## kra\_creatp Subroutine

### Purpose

Creates a new kernel process and attaches it to a resource set.

### Syntax

```
#include <sys/rset.h>
int kra_creatp (pid, rstype, rsid, flags)
pid_t *pid;
rstype_t rstype;
rsid_t rsid;
unsigned int flags;
```

### Description

The **kra\_creatp** kernel service creates a new kernel process and attaches it to a resource set. The **kra\_creatp** kernel service attaches the new kernel process to the resource set specified by the *rstype* and *rsid* parameters.

The **kra\_creatp** kernel service is similar to the **creatp** kernel service.

The following conditions must be met to successfully attach a kernel process to a resource set:

- The resource set must contain processors that are available in the system.
- The calling process must either have root authority or have CAP\_NUMA\_ATTACH capability.
- The calling thread must not have a bindprocessor binding to a processor.
- The resource set must be contained in the calling process' partition resource set.

**Note:** When the **creatp** kernel service is used, the new kernel process inherits its parent's resource set attachments.

### Parameters

Item	Description
<i>pid</i>	Pointer to a <b>pid_t</b> field to receive the process ID of the new kernel process.
<i>rstype</i>	Specifies the type of resource the new process will be attached to. This parameter must be the following value, defined in <b>rset.h</b> . <ul style="list-style-type: none"><li>• R_RSET: resource set.</li></ul>
<i>rsid</i>	Identifies the resource set the new process will be attached to. <ul style="list-style-type: none"><li>• Resource set ID (for <i>rstype</i> of R_RSET): set the <i>rsid_t at_rset</i> field to the desired resource set.</li></ul>
<i>flags</i>	Reserved for future use. Specify as 0.

### Return Values

Upon successful completion, the **kra\_creatp** kernel service returns a 0. If unsuccessful, one or more of the following are true:

Item	Description
EINVAL	One of the following is true: <ul style="list-style-type: none"> <li>• The <i>rstype</i> parameter contains an invalid type identifier.</li> <li>• The <i>flags</i> parameter contains an invalid flags value.</li> </ul>
ENODEV	The specified resource set does not contain any available processors.
EFAULT	Invalid address.
EPERM	One of the following is true: <ul style="list-style-type: none"> <li>• The calling process has neither root authority nor CAP_NUMA_ATTACH attachment privilege.</li> <li>• The calling process contains one or more threads with a bindprocessor processor binding.</li> <li>• The specified resource set is not included in the calling process' partition resource set.</li> </ul>
ENOMEM	Memory not available.

#### Related reference:

“creatp Kernel Service” on page 54

“initp Kernel Service” on page 221

“kra\_attachrset Subroutine” on page 293

## kra\_detachrset Subroutine

### Purpose

Detaches a work component from a resource set.

### Syntax

```
#include <sys/rset.h>
int kra_detachrset (rstype, rsid, flags)
rstype_t rstype;
rsid_t rsid;
unsigned int flags;
```

### Description

The **kra\_detachrset** subroutine detaches a work component specified by *rstype* and *rsid* from a resource set.

The work component is an existing process identified by the process ID or an existing kernel thread identified by the kernel thread ID (tid). A process ID or thread ID value of RS\_MYSELF indicates the detach command applies to the current process or the current kernel thread, respectively.

The following conditions must be met to detach a process or kernel thread from a resource set:

- The calling process must either have root authority or have CAP\_NUMA\_ATTACH capability.
- The calling process must either have root authority or the same effective userid as the target process.

If these conditions are not met, the operation will fail.

Once a process is detached from a resource set, the threads in the process can run on all available processors contained in the process' partition resource set. Once a kernel thread is detached from a resource set, that thread can run on all available processors contained in its process effective or partition resource set.

### Parameters



Item	Description
<i>rstype</i>	Specifies the type of work component to be detached from to the resource set specified by <i>rset</i> . This parameter must be the following value, defined in <b>rset.h</b> : <ul style="list-style-type: none"> <li>R_PROCESS: existing process</li> <li>R_THREAD: existing kernel thread</li> </ul>
<i>rsid</i>	Identifies the work component to be attached to the resource set specified by <i>rset</i> . This parameter must be the following: <ul style="list-style-type: none"> <li>Process ID (for <i>rstype</i> of R_PROCESS): set the <i>rsid_t at_pid</i> field to the desired process' process ID.</li> <li>Kernel thread ID (for <i>rstype</i> of R_THREAD): set the <i>rsid_t at_tid</i> field to the desired kernel thread's thread ID.</li> </ul>
<i>flags</i>	For <i>rstype</i> of R_PROCESS, the R_DETACH_ALLTHRDS indicates that R_THREAD <i>rsets</i> are detached from all threads in a specified process. The process' effective <i>rset</i> is not detached in this case. Reserved for future use. Specify as 0.

## Return Values

Upon successful completion, the **kra\_detachrset** subroutine returns a 0. If unsuccessful, one or more of the following are true:

Item	Description
EINVAL	One of the following is true: <ul style="list-style-type: none"> <li>The <i>flags</i> parameter contains an invalid value.</li> <li>The <i>rstype</i> contains an invalid type qualifier.</li> </ul>
ESRCH	The process or kernel thread identified by the <i>rstype</i> and <i>rsid</i> parameters does not exist.
EPERM	One of the following is true: <ul style="list-style-type: none"> <li>The calling process has neither root authority nor CAP_NUMA_ATTACH attachment privilege.</li> <li>The calling process has neither root authority nor the same effective user ID as the process identified by the <i>rstype</i> and <i>rsid</i> parameters.</li> </ul>

### Related reference:

“kra\_attachrset Subroutine” on page 293

## kra\_getrset Subroutine

### Purpose

Gets the resource set to which a work component is attached.

### Syntax

```
#include <sys/rset.h>
int kra_getrset (rstype, rsid, flags, rset, rset_type)
rstype_t rstype;
rsid_t rsid;
unsigned int flags;
rsethandle_t rset;
unsigned int *rset_type;
```

### Description

The **kra\_getrset** subroutine returns the resource set to which a specified work component is attached.

The work component is an existing process identified by the process ID or an existing kernel thread identified by the kernel thread ID (tid). A process ID or thread ID value of RS\_MYSELF indicates the resource set attached to the current process or the current kernel thread, respectively, is requested.

Upon successful completion, one of the following types of resource set is returned into the *rset\_type* parameter:

- A value of RS\_EFFECTIVE\_RSET indicates the process was explicitly attached to the resource set. This may have been done with the **kra\_attachrset** subroutine.

- A value of RS\_PARTITION\_RSET indicates the process was not explicitly attached to a resource set. However, the process had an explicitly set partition resource set. This may be set with the `krs_setpartition` subroutine or through the use of WLM work classes with resource sets.
- A value of RS\_DEFAULT\_RSET indicates the process was not explicitly attached to a resource set nor did it have an explicitly set partition resource set. The system default resource set is returned.
- A value of RS\_THREAD\_RSET indicates the kernel thread was explicitly attached to the resource set. This might have been done with the `ra_attachrset` subroutine.

## Parameters

Item	Description
<code>rstype</code>	Specifies the type of the work component whose resource set attachment is requested. This parameter must be the following value, defined in <code>rset.h</code> : <ul style="list-style-type: none"> <li>• R_PROCESS: existing process</li> <li>• R_THREAD: existing kernel thread</li> </ul>
<code>rsid</code>	Identifies the work component whose resource set attachment is requested. This parameter must be the following: <ul style="list-style-type: none"> <li>• Process ID (for <code>rstype</code> of R_PROCESS): set the <code>rsid_t at_pid</code> field to the desired process' process ID.</li> <li>• Kernel thread ID (for <code>rstype</code> of R_THREAD): set the <code>rsid_t at_tid</code> field to the desired kernel thread's thread ID.</li> </ul>
<code>flags</code>	Reserved for future use. Specify as 0.
<code>rset</code>	Specifies the resource set to receive the work component's resource set.
<code>rset_type</code>	Points to an unsigned integer field to receive the resource set type.

## Return Values

Upon successful completion, the `kra_getrset` subroutine returns a 0. If unsuccessful, one or more of the following are true:

Item	Description
EINVAL	One of the following is true: <ul style="list-style-type: none"> <li>• The <code>flags</code> parameter contains an invalid value.</li> <li>• The <code>rstype</code> parameter contains an invalid type qualifier.</li> </ul>
EFAULT	Invalid address.
ESRCH	The process or kernel thread identified by the <code>rstype</code> and <code>rsid</code> parameters does not exist.

### Related reference:

“`krs_getpartition` Subroutine” on page 302

## `krs_alloc` Subroutine

### Purpose

Allocates a resource set and returns its handle.

### Syntax

```
#include <sys/rset.h>
int krs_alloc (rset, flags)
rsethandle_t *rset;
unsigned int flags;
```

### Description

The `krs_alloc` subroutine allocates a resource set and initializes it according to the information specified by the `flags` parameter. The value of the `flags` parameter determines how the new resource set is initialized.

### Parameters

Item	Description
<i>rset</i>	Points to an <code>rsethandle_t</code> where the resource set handle is stored on successful completion.
<i>flags</i>	Specifies how the new resource set is initialized. It takes one of the following values, defined in <code>rset.h</code> : <ul style="list-style-type: none"> <li>• <b>RS_EMPTY</b> (or 0 value): The resource set is initialized to contain no resources.</li> <li>• <b>RS_SYSTEM</b>: The resource set is initialized to contain available system resources.</li> <li>• <b>RS_ALL</b>: The resource set is initialized to contain all resources.</li> <li>• <b>RS_PARTITION</b>: The resource set is initialized to contain the resources in the caller's process partition resource set.</li> </ul>

## Return Values

Upon successful completion, the `krs_alloc` subroutine returns a 0. If unsuccessful, one or more of the following is returned:

Item	Description
EINVAL	The <i>flags</i> parameter contains an invalid value.
ENOMEM	There is not enough space to create the data structures related to the resource set.

### Related reference:

“`krs_free` Subroutine”

“`krs_getinfo` Subroutine” on page 301

“`krs_init` Subroutine” on page 304

## `krs_free` Subroutine

### Purpose

Frees a resource set.

### Syntax

```
#include <sys/rset.h>
void krs_free(rset)
rsethandle_t rset;
```

### Description

The `krs_free` subroutine frees a resource set identified by the *rset* parameter. The resource set must have been allocated by the `krs_alloc` subroutine.

### Parameters

Item	Description
<i>rset</i>	Specifies the resource set whose memory will be freed.

### Related reference:

“`krs_alloc` Subroutine” on page 298

## `krs_getassociativity` Subroutine

### Purpose

Gets the hardware associativity values for a resource.

## Syntax

```
#include <sys/rset.h>
int krs_getassociativity (type, id, assoc_array, array_size)
unsigned int type;
unsigned int id;
unsigned int *assoc_array;
unsigned int array_size;
```

## Description

The **krs\_getassociativity** subroutine returns the array of hardware associativity values for a specified resource.

This is a special purpose subroutine intended for specialized root applications needing the hardware associativity value information. The **krs\_getinfo**, **krs\_getrad**, and **krs\_numrads** subroutines are provided for typical applications to discover system hardware topology.

The calling process must have root authority to get hardware associativity values.

## Parameters

Item	Description
<i>type</i>	Specifies the resource type whose associativity values are requested. The only value supported to retrieve values for a processor is R_PROCS.
<i>id</i>	Specifies the logical resource id whose associativity values are requested.
<i>assoc_array</i>	Specifies the address of an array of unsigned integers to receive the associativity values.
<i>array_size</i>	Specifies the number of unsigned integers in <i>assoc_array</i> .

## Return Values

Upon successful completion, the **krs\_getassociativity** subroutine returns a 0. The *assoc\_array* parameter array contains the resource's associativity values. The first entry in the array indicates the number of associativity values returned. If the hardware system does not provide system topology data, a value of 0 is returned in the first array entry. If unsuccessful, one or more of the following are returned:

Item	Description
EINVAL	One of the following occurred: <ul style="list-style-type: none"><li>The <i>array_size</i> parameter was specified as 0.</li><li>An invalid <i>type</i> parameter was specified.</li></ul>
ENODEV	The resource specified by the <i>id</i> parameter does not exist.
EFAULT	Invalid address.
EPERM	The calling process does not have root authority.

### Related reference:

“krs\_getinfo Subroutine” on page 301

“krs\_getrad Subroutine” on page 303

“krs\_numrads Subroutine” on page 305

## krs\_get\_homesrad Subroutine

### Purpose

Gets the currently running thread's home SRADID (Scheduler Resource Allocation Domain Identifier).

### Library

Standard C library (**libc.a**)

## Syntax

```
#include <sys/rset.h>
sradid_t krs_get_homesrad(void)
```

## Description

The `krs_get_homesrad` is a kernel service and if the `ENHANCED_AFFINITY` services are enabled, the `krs_get_homesrad` subroutine returns the home SRADID of the currently running thread. If the `ENHANCED_AFFINITY` services are not enabled, the `krs_get_homesrad` subroutine returns `SRADID_ANY`. SRADID is the index of a RAD (Resource Allocation Domain) at the `R_SRADSDL` system detail level.

## Return Values

If the `ENHANCED_AFFINITY` services are enabled, the home SRADID of the currently running thread is returned. Otherwise, `SRADID_ANY` is returned.

### Related reference:

“`krs_getrad` Subroutine” on page 303

## `krs_getinfo` Subroutine

### Purpose

Gets information about a resource set.

## Syntax

```
#include <sys/rset.h>
int krs_getinfo(rset, info_type, flags, result)
rsethandle_t rset;
rsinfo_t info_type;
unsigned int flags;
int *result;
```

## Description

The `krs_getinfo` subroutine retrieves information about the resource set identified by the `rset` parameter. Depending on the value of the `info_type` parameter, the `krs_getinfo` subroutine returns information about the number of available processors, the number of available memory pools, or the amount of available memory contained in the resource `rset`.

The subroutine can also return global system information such as the maximum system detail level, the symmetric multiprocessor (SMP) and multiple chip module (MCM) system detail levels, and the maximum number of processor or memory pool resources in a resource set.

## Parameters

Item	Description
<code>rset</code>	Specifies a resource set handle of a resource set the information should be retrieved from. This parameter is not meaningful if the <code>info_type</code> parameter is <code>R_MAXSDL</code> , <code>R_MAXPROCS</code> , <code>R_MAXMEMPS</code> , <code>R_SMPSDL</code> , or <code>R_MCMSDL</code> .

Item	Description
<i>info_type</i>	Specifies the type of information being requested. One of the following values (defined in <code>rset.h</code> ) can be used: <ul style="list-style-type: none"> <li>• <b>R_NUMPROCS</b>: The number of available processors in the resource set is returned.</li> <li>• <b>R_NUMMEMPS</b>: The number of available memory pools in the resource set is returned.</li> <li>• <b>R_MEMSIZE</b>: The amount of available memory (in MB) contained in the resource set is returned.</li> <li>• <b>R_MAXSDL</b>: The maximum system detail level of the system is returned.</li> <li>• <b>R_MAXPROCS</b>: The maximum number of processors that may be contained in a resource set is returned.</li> <li>• <b>R_MAXMEMPS</b>: The maximum number of memory pools that may be contained in a resource set is returned.</li> <li>• <b>R_SMPSDL</b>: The system detail level that corresponds to the traditional notion of an SMP is returned. A system detail level of 0 is returned if the hardware system does not provide system topology data.</li> <li>• <b>R_MCMSDL</b>: The system detail level that corresponds to resources packaged in an MCM is returned. A system detail level of 0 is returned if the hardware system does not have MCMs or does not provide system topology data.</li> </ul>
<i>flags</i>	Reserved for future use. Must be specified as 0.
<i>result</i>	Points to an integer where the result is stored on successful completion.

## Return Values

Upon successful completion, the `krs_getinfo` subroutine returns a 0, and the *result* field contains the requested information. If unsuccessful, one or more of the following are returned:

Item	Description
EINVAL	One of the following is true: <ul style="list-style-type: none"> <li>• The <i>info_type</i> parameter specifies an invalid resource type value.</li> <li>• The <i>flags</i> parameter was not specified as 0.</li> </ul>
EFAULT	Invalid address.

### Related reference:

“`krs_numrads` Subroutine” on page 305

## `krs_getpartition` Subroutine Purpose

Gets the partition resource set to which a process is attached.

### Syntax

```
#include <sys/rset.h>
int krs_getpartition (pid, flags, rset, rset_type)
pid_t pid;
unsigned int flags;
rsethandle_t rset;
unsigned int *rset_type;
```

### Description

The `krs_getpartition` subroutine returns the partition resource set attached to the specified process. A process ID value of `RS_MYSELF` indicates the partition resource set attached to the current process is requested.

Upon successful completion, the type of resource set is returned into the *rset\_type* parameter.

A value of `RS_PARTITION_RSET` indicates the process has a partition resource set that is set explicitly. This may be set with the `krs_setpartition` subroutine or through the use of WLM work classes with resource sets.

A value of `RS_DEFAULT_RSET` indicates the process did not have an explicitly set partition resource set. The system default resource set is returned.

## Parameters

Item	Description
<i>pid</i>	Specifies the process ID whose partition <i>rset</i> is requested.
<i>flags</i>	Reserved for future use. Specify as 0.
<i>rset</i>	Specifies the resource set to receive the process' partition resource set.
<i>rset_type</i>	Points to an unsigned integer field to receive the resource set type.

## Return Values

Upon successful completion, the `krs_getpartition` subroutine returns a 0. If unsuccessful, one or more of the following are true:

Item	Description
EFAULT	Invalid address.
ESRCH	The process identified by the <i>pid</i> parameter does not exist.

### Related reference:

“`kra_getrset` Subroutine” on page 297

## `krs_getrad` Subroutine Purpose

Returns a system resource allocation domain (RAD) contained in an input resource set.

### Syntax

```
#include <sys/rset.h>
int krs_getrad (rad, sdl, index, flags)
rsethandle_t rad;
unsigned int sdl;
unsigned int index;
unsigned int flags;
```

### Description

The `krs_getrad` subroutine returns a system RAD at a specified system detail level and index.

The system RAD is specified by system detail level *sdl* and index number *index*.

The *rad* parameter must be allocated (using the `krs_alloc` subroutine) prior to calling the `krs_getrad` subroutine.

## Parameters

Item	Description
<i>rad</i>	Specifies a resource set handle to receive the desired system RAD.
<i>sdl</i>	Specifies the system detail level of the desired system RAD.
<i>index</i>	Specifies the index of the system RAD that should be returned from among those at the specified <i>sdl</i> . This parameter must belong to the <code>[0, krs_numrads(rset, sdl, flags) - 1]</code> interval.
<i>flags</i>	Reserved for future use. Specify as 0.

## Return Values

Upon successful completion, the `krs_getrad` subroutine returns a 0. If unsuccessful, one or more of the following are true:

Item	Description
EINVAL	One of the following is true: <ul style="list-style-type: none"><li>The <i>flags</i> parameter contains an invalid value.</li><li>The <i>sdl</i> parameter is greater than the maximum system detail level.</li><li>The RAD specified by the <i>index</i> parameter does not exist at the system detail level specified by the <i>sdl</i> parameter.</li></ul>
EFAULT	Invalid address.

### Related reference:

“`krs_numrads` Subroutine” on page 305

“`krs_getinfo` Subroutine” on page 301

“`krs_alloc` Subroutine” on page 298

## `krs_init` Subroutine

### Purpose

Initializes a previously allocated resource set.

### Syntax

```
#include <sys/rset.h>
int krs_init (rset, flags)
rsethandle_t rset;
unsigned int flags;
```

### Description

The `krs_init` subroutine initializes a previously allocated resource set. The resource set is initialized according to information specified by the *flags* parameter.

### Parameters

Item	Description
<i>rset</i>	Specifies the handle of the resource set to initialize.
<i>flags</i>	Specifies how the resource set is initialized. It takes one of the following values, defined in <code>rset.h</code> : <ul style="list-style-type: none"><li><b>RS_EMPTY</b>: The resource set is initialized to contain no resources.</li><li><b>RS_SYSTEM</b>: The resource set is initialized to contain available system resources.</li><li><b>RS_ALL</b>: The resource set is initialized to contain all resources.</li><li><b>RS_PARTITION</b>: The resource set is initialized to contain the resources in the caller's process partition resource set.</li></ul>

## Return Values

Upon successful completion, the `krs_init` subroutine returns a 0. If unsuccessful, the following is returned:



Item	Description
EINVAL	The <i>flags</i> parameter contains an invalid value.

**Related reference:**

“krs\_alloc Subroutine” on page 298

## krs\_numrads Subroutine

### Purpose

Returns the number of system resource allocation domains (RADs) that have available resources.

### Syntax

```
#include <sys/rset.h>
int krs_numrads(rset, sdl, flags)
rsethandle_t rset;
unsigned int sdl;
unsigned int flags;
```

### Description

The **krs\_numrads** subroutine returns the number of system RADs at system detail level *sdl*, that have available resources contained in the resource set identified by the *rset* parameter.

The number of atomic RADs contained in the *rset* parameter is returned if the *sdl* parameter is equal to the maximum system detail level.

### Parameters

Item	Description
<i>rset</i>	Specifies the resource set handle for the resource set being queried.
<i>sdl</i>	Specifies the system detail level in which the caller is interested.
<i>flags</i>	Reserved for future use. Specify as 0.

### Return Values

Upon successful completion, the number of RADs is returned. If unsuccessful, a -1 is returned and one or more of the following are true:

- The *flags* parameter contains an invalid value.
- The *sdl* parameter is greater than the maximum system detail level.

**Related reference:**

“krs\_getrad Subroutine” on page 303

“krs\_getinfo Subroutine” on page 301

## krs\_op Subroutine

### Purpose

Performs a set of operations on one or two resource sets.

### Syntax

```
#include <sys/rset.h>
int krs_op(command, rset1, rset2, flags, id)
unsigned int command;
rsethandle_t rset1, rset2;
unsigned int flags;
unsigned int id;
```

## Description

The `kr_s_op` subroutine performs the operation specified by the `command` parameter on resource set `rset1`, or both resource sets `rset1` and `rset2`.

## Parameters

Item	Description
<code>command</code>	<p>Specifies the operation to apply to the resource sets identified by <code>rset1</code> and <code>rset2</code>. One of the following values, defined in <code>rset.h</code>, can be used:</p> <ul style="list-style-type: none"><li>• <b>RS_UNION</b>: The resources contained in either <code>rset1</code> or <code>rset2</code> are stored in <code>rset2</code>.</li><li>• <b>RS_INTERSECTION</b>: The resources that are contained in both <code>rset1</code> and <code>rset2</code> are stored in <code>rset2</code>.</li><li>• <b>RS_EXCLUSION</b>: The resources in <code>rset1</code> that are also in <code>rset2</code> are removed from <code>rset2</code>. On completion, <code>rset2</code> contains all the resources that were contained in <code>rset2</code> but were not contained in <code>rset1</code>.</li><li>• <b>RS_COPY</b>: All resources in <code>rset1</code> whose type is <code>flags</code> are stored in <code>rset2</code>. If <code>rset1</code> contains no resources of this type, <code>rset2</code> will be empty. The previous content of <code>rset2</code> is lost, while the content of <code>rset1</code> is unchanged.</li><li>• <b>RS_IEMPTY</b>: Test if resource set <code>rset1</code> is empty.</li><li>• <b>RS_ISEQUAL</b>: Test if resource sets <code>rset1</code> and <code>rset2</code> are equal.</li><li>• <b>RS_ISCONTAINED</b>: Test if all resources in resource set <code>rset1</code> are also contained in resource set <code>rset2</code>.</li><li>• <b>RS_TESTRESOURCE</b>: Test if the resource whose type is <code>flags</code> and index is <code>id</code> is contained in resource set <code>rset1</code>.</li><li>• <b>RS_ADDRESOURCE</b>: Add the resource whose type is <code>flags</code> and index is <code>id</code> to resource set <code>rset1</code>.</li><li>• <b>RS_DELRESOURCE</b>: Delete the resource whose type is <code>flags</code> and index is <code>id</code> from resource set <code>rset1</code>.</li><li>• <b>RS_STSET</b>: Constructs an ST resource set by including only one hardware thread per physical processor included in <code>rset1</code> and stores it in <code>rset2</code>. Only available processors are considered when constructing the ST resource set.</li></ul>
<code>rset1</code>	Specifies the resource set handle for the first of the resource sets involved in the <code>command</code> operation.
<code>rset2</code>	Specifies the resource set handle for the second of the resource sets involved in the <code>command</code> operation. This resource set is also used, on return, to store the result of the operation, and its previous content is lost. The <code>rset2</code> parameter is ignored on the <code>RS_IEMPTY</code> , <code>RS_TESTRESOURCE</code> , <code>RS_ADDRESOURCE</code> , and <code>RS_DELRESOURCE</code> commands.
<code>flags</code>	<p>When combined with the <code>RS_COPY</code> command, the <code>flags</code> parameter specifies the type of the resources that will be copied from <code>rset1</code> to <code>rset2</code>. This parameter is constructed by logically ORing one or more of the following values, defined in <code>rset.h</code>:</p> <ul style="list-style-type: none"><li>• <b>R_PROCS</b>: processors</li><li>• <b>R_MEMPS</b>: memory pools</li><li>• <b>R_ALL_RESOURCES</b>: processors and memory pools</li></ul> <p>If none of the above are specified for <code>flags</code>, <code>R_ALL_RESOURCES</code> is assumed.</p>
<code>id</code>	<p>On the <code>RS_TESTRESOURCE</code>, <code>RS_ADDRESOURCE</code>, and <code>RS_DELRESOURCE</code> commands, the <code>id</code> parameter specifies the index of the resource to be tested, added, or deleted. This parameter is ignored on the other commands.</p>

## Return Values

Item	Description
0	Successful completion. The tested condition is not met for the <code>RS_IEMPTY</code> , <code>RS_ISEQUAL</code> , <code>RS_ISCONTAINED</code> , and <code>RS_TESTRESOURCE</code> commands.
1	Successful completion. The tested condition is met for the <code>RS_IEMPTY</code> , <code>RS_ISEQUAL</code> , <code>RS_ISCONTAINED</code> , and <code>RS_TESTRESOURCE</code> commands.
-1	<p>Unsuccessful completion. One or more of the following are true:</p> <ul style="list-style-type: none"><li>• <code>rset1</code> identifies an invalid resource set.</li><li>• <code>rset2</code> identifies an invalid resource set.</li><li>• <code>command</code> identifies an invalid operation.</li><li>• <code>flags</code> identifies an invalid resource type.</li><li>• <code>id</code> specifies a resource index that is too large.</li><li>• Invalid address.</li></ul>

## krs\_setpartition Subroutine

### Purpose

Sets the partition resource set of a process.

### Syntax

```
#include <sys/rset.h>
int krs_setpartition(pid, rset, flags)
pid_t pid;
rsethandle_t rset;
unsigned int flags;
```

### Description

The **krs\_setpartition** subroutine sets a process' partition resource set. The subroutine can also be used to remove a process' partition resource set.

The partition resource set limits the threads in a process to running only on the processors contained in the partition resource set.

The work component is an existing process identified by process ID. A process ID value of RS\_MYSELF indicates the attachment applies to the current process.

The following conditions must be met to set a process' partition resource set:

- The calling process must have root authority.
- The resource set must contain processors that are available in the system.
- The new partition resource set must be equal to, or a superset of the target process' effective resource set.
- The target process must not contain any threads that have bindprocessor bindings to a processor.

The *flags* parameter can be set to indicate the policy for using the resources contained in the resource set specified in the *rset* parameter. The only supported scheduling policy is **R\_ATTACH\_STRSET**, which is useful only when the processors of the system are running in simultaneous multithreading mode. Processors like the POWER5 support simultaneous multithreading, where each physical processor has two execution engines, called *hardware threads*. Each hardware thread is essentially equivalent to a single CPU, and each is identified as a separate CPU in a resource set. The **R\_ATTACH\_STRSET** flag indicates that the process is to be scheduled with a single-threaded policy; namely, that it should be scheduled on only one hardware thread per physical processor. If this flag is specified, then all of the available processors indicated in the resource set must be of exclusive use. A new resource set, called an *ST resource set*, is constructed from the specified resource set and attached to the process according to the following rules:

- All offline processors are ignored.
- If all the hardware threads (CPUs) of a physical processor (when running in simultaneous multithreading mode, there will be more than one active hardware thread per physical processor) are not included in the specified resource set, the other CPUs of the processor are ignored when constructing the ST resource set.
- Only one CPU (hardware thread) resource per physical processor is included in the ST resource set.

### Parameters

Item	Description
<i>pid</i>	Specifies the process ID of the process whose partition resource set is to be set. A value of RS_MYSELF indicates the current process' partition resource set should be set.
<i>rset</i>	Specifies the partition resource set to be set. A value of RS_DEFAULT indicates the process' partition resource set should be removed.
<i>flags</i>	Specifies the policy to use for the process. A value of R_ATTACH_STRSET indicates that the process is to be scheduled with a single-threaded policy (only on one hardware thread per physical processor).

## Return Values

Upon successful completion, the **krs\_setpartition** subroutine returns a 0. If unsuccessful, one or more of the following are true:

Item	Description
EINVAL	The R_ATTACH_STRSET <i>flags</i> parameter is specified and one or more processors in the <i>rset</i> parameter are not assigned for exclusive use.
ENODEV	The resource set specified by the <i>rset</i> parameter does not contain any available processors, or the R_ATTACH_STRSET <i>flags</i> parameter is specified and the constructed ST resource set does not have any available processors.
ESRCH	The process identified by the <i>pid</i> parameter does not exist.
EFAULT	Invalid address.
ENOMEM	Memory not available.
EPERM	One of the following is true: <ul style="list-style-type: none"> <li>• The calling process does not have root authority.</li> <li>• The process identified by the <i>pid</i> parameter has one or more threads with a bindprocessor processor binding.</li> <li>• The process identified by the <i>pid</i> parameter has an effective resource set and the new partition resource set identified by the <i>rset</i> parameter does not contain all of the effective resource set's resources.</li> </ul>

### Related reference:

“krs\_getpartition Subroutine” on page 302

“kra\_attachrset Subroutine” on page 293

### Related information:

Exclusive use processor resource sets

## ksettckd Kernel Service Purpose

Sets the current status of the systemwide timer-adjustment values.

### Syntax

```
#include <sys/types.h>
int ksettckd (timed, tickd, time_adjusted)
int *timed;
int *tickd;
int *time_adjusted;
```

### Parameters

Item	Description
<i>timed</i>	Specifies the number of microseconds by which the systemwide timer is to be adjusted unless set to a null pointer.
<i>tickd</i>	Specifies the adjustment rate of the systemwide timer unless set to a null pointer. This rate determines the number of microseconds that the systemwide timer is adjusted with each timer tick. Adjustment continues until the time has been corrected by the amount specified by the <i>timed</i> parameter.
<i>time_adjusted</i>	Sets the kernel-maintained time adjusted flag to True or False. If the <i>time_adjusted</i> parameter is a null pointer, calling the <b>ksettickd</b> kernel service always sets the kernel's <i>time_adjusted</i> parameter to False.

## Description

The **ksettickd** kernel service provides kernel extensions with the capability to update the *time\_adjusted* parameter, and set or change the systemwide time-of-day timer adjustment amount and rate. The timer-adjustment values indicated by the *timed* and *tickd* parameters are the same values used by the **adjtime** subroutine. A call to the **settimer** or **adjtime** subroutine for the systemwide time-of-day timer sets the *time\_adjusted* parameter to True, as read by the **kgettickd** kernel service.

This kernel service is typically used only by kernel extensions providing time synchronization functions such as coordinated network time where the **adjtime** subroutine is insufficient.

**Note:** The **ksettickd** service provides no serialization with respect to the **adjtime** and **settimer** subroutines, the **ksettimer** kernel service, or the timer interrupt handler, all of which also use and update these values. The caller of this kernel service must provide the necessary serialization to ensure appropriate operation.

## Execution Environment

The **ksettickd** kernel service can be called from either the process or interrupt environment.

## Return Value

The **ksettickd** kernel service always returns a value of 0.

### Related reference:

“kgettickd Kernel Service” on page 271

### Related information:

adjtime subroutine

Using Fine Granularity Timer Services and Structures

## ksettimer Kernel Service

### Purpose

Sets the systemwide time-of-day timer.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/time.h>

int ksettimer (nct)
struct timestruc_t *nct;
```

### Parameter

Item	Description
------	-------------

<i>nct</i>	Points to a <code>timestruc_t</code> structure, which contains the new current time to be set. The nanoseconds member of this structure is valid only if greater than or equal to 0, and less than the number of nanoseconds in a second.
------------	---

## Description

The `ksettimer` kernel service provides a kernel extension with the capability to set the systemwide time-of-day timer. Kernel extensions typically use this kernel service to support network coordinated time, which is the periodic synchronization of all system clocks to a common time by a time server or set of time servers on a network. The newly set "current" time must represent the amount of time since 00:00:00 GMT, January 1, 1970.

## Execution Environment

The `ksettimer` kernel service can be called from the process environment only.

## Return Values

Item	Description
------	-------------

0	Indicates success.
---	--------------------

EINVAL	Indicates that the new current time specified by the <i>nct</i> parameter is outside the range of the systemwide timer.
--------	---

EIO	Indicates that an error occurred while this kernel service was accessing the timer device.
-----	--

## Related information:

Using Fine Granularity Timer Services and Structures

Timer and Time-of-Day Kernel Services

## kthread\_kill Kernel Service

### Purpose

Posts a signal to a specified kernel-only thread.

### Syntax

```
#include <sys/thread.h>
```

```
void kthread_kill ( tid, sig)
```

```
tid_t tid;
```

```
int sig;
```

### Parameters

Item	Description
------	-------------

<i>tid</i>	Specifies the target kernel-only thread. If its value is -1, the signal is posted to the calling thread.
------------	--

<i>sig</i>	Specifies the signal number to post.
------------	--------------------------------------

## Description

The `kthread_kill` kernel service posts the signal *sig* to the kernel thread specified by the *tid* parameter. When the service is called from the process environment, the target thread must be in the same process as the calling thread. When the service is called from the interrupt environment, the signal is posted to the target thread, without a permission check.

## Execution Environment

The `kthread_kill` kernel service can be called from either the process environment or the interrupt environment.

## Return Values

The `kthread_kill` kernel service has no return values.

### Related reference:

“`sig_chk` Kernel Service” on page 473

### Related information:

Process and Exception Management Kernel Services

## `kthread_start` Kernel Service

### Purpose

Starts a previously created kernel-only thread.

### Syntax

```
#include <sys/thread.h>
```

```
int kthread_start ( tid, i_func, i_data_addr, i_data_len, i_stackaddr,
i_sigmask)
tid_t tid;
int (*i_func) (void *);
void *i_data_addr;
size_t i_data_len;
void *i_stackaddr;
sigset_t *i_sigmask;
```

### Parameters

Item	Description
<i>tid</i>	Specifies the kernel-only thread to start.
<i>i_func</i>	Points to the entry-point routine of the kernel-only thread.
<i>i_data_addr</i>	Points to data that will be passed to the entry-point routine.
<i>i_data_len</i>	Specifies the length of the data chunk.
<i>i_stackaddr</i>	Specifies the stack's base address for the kernel-only thread.
<i>i_sigmask</i>	Specifies the set of signal to block from delivery when the new kernel-only thread begins execution.

### Description

The `kthread_start` kernel service starts the kernel-only thread specified by the *tid* parameter. The thread must have been previously created with the `thread_create` kernel service, and its state must be **TSIDL**.

This kernel service initializes and schedules the thread for the processor. Its state is changed to **TSRUN**. The thread is initialized so that it begins executing at the entry point specified by the *i\_func* parameter, and that the signals specified by the *i\_sigmask* parameter are blocked from delivery.

The thread's entry point gets one parameter, a pointer to a chunk of data that is copied to the base of the thread's stack. The *i\_data\_addr* and *i\_data\_len* parameters specify the location and quantity of data to copy. The format of the data must be agreed upon by the initializing and initialized thread.

The thread's stack's base address is specified by the *i\_stackaddr* parameter. If a value of zero is specified, the kernel will allocate the memory for the stack (96K). This memory will be reclaimed by the system

when the thread terminates. If a non-zero value is specified, then the caller should allocate the backing memory for the stack. Since stacks grow from high addresses to lower addresses, the *i\_stackaddr* parameter specifies the highest address for the thread's stack.

The thread will be automatically terminated when it returns from the entry point routine. If it is the last thread in the process, then the process will be exited.

## Execution Environment

The `kthread_start` kernel service can be called from the process environment only.

## Return Values

The `kthread_start` kernel service returns one of the following values:

Item	Description
0	Indicates a successful start.
ESRCH	Indicates that the <i>tid</i> parameter is not valid.

### Related reference:

“thread\_create Kernel Service” on page 485

### Related information:

Process and Exception Management Kernel Services

## kvmgetinfo Kernel Service

### Purpose

Retrieves Virtual Memory Manager (VMM) information.

### Syntax

```
#include <sys/vminfo.h>
```

```
int kvmgetinfo ( void *out, int command, int arg)
```

### Description

The `kvmgetinfo` kernel service returns the current value of certain VMM parameters.

### Parameters

Item	Description
<i>out</i>	Specifies the address where VMM information should be returned.



Item	Description
<i>command</i>	<p>Specifies which information should be returned. The valid values for the <i>command</i> parameter are described below:</p> <p><b>VMINFO</b> The content of <b>vminfo</b> structure (described in <b>sys/vminfo.h</b>) will be returned. The <i>out</i> parameter should point to a <b>vminfo</b> structure and the <i>arg</i> parameter should be the size of this structure. The smaller of the <i>arg</i> or <i>sizeof</i> (<b>struct vminfo</b>) parameters will be copied.</p> <p><b>VMINFO_ABRIDGED</b> The content of the <b>vminfo</b> structure (described in the <b>sys/vminfo.h</b> file) is returned. For this command, only the non-time consuming statistics are updated, so this command must be used in performance-critical applications rather than the <b>VMINFO</b> command. The <i>out</i> parameter must point to a <b>vminfo</b> structure and the <i>arg</i> parameter must be the size of this structure. The smaller of the <i>arg</i> or <i>sizeof</i> (<b>struct vminfo</b>) parameters are copied.</p> <p><b>VM_PAGE_INFO</b> The size, in bytes, of the page backing the address specified in the <i>addr</i> field of the <b>vm_page_info</b> structure (described in the <b>sys/vminfo.h</b> file) is returned. The <i>out</i> parameter should point to a <b>vm_page_info</b> structure with the <i>addr</i> field set to the desired address of which to query the page size. This address, <i>addr</i>, is interpreted as an address in the address space of the current running process. The <i>arg</i> parameter should be the size of the <b>vm_page_info</b> structure.</p> <p><b>IPC_LIMITS</b> The content of the <b>ipc_limits</b> struct (described in the <b>sys/vminfo.h</b> file) is returned. The <i>out</i> parameter should point to an <b>ipc_limits</b> structure and <i>arg</i> should be the size of this structure. The smaller of the <i>arg</i> or <i>sizeof</i> (<b>struct ipc_limits</b>) parameters will be copied. The <b>ipc_limits</b> struct contains the inter-process communication (IPC) limits for the system.</p> <p><b>VMINFO_GETPSIZES</b> Reports a system's supported page sizes. When <i>arg</i> is 0, the <i>out</i> parameter is ignored, and the number of supported page sizes is returned. When <i>arg</i> is greater than 0, <i>arg</i> indicates the number of page sizes to report, and <i>out</i> must be a pointer to an array with <i>arg</i> number of <b>psize_t</b> types. The array of <b>psize_t</b> types is updated with the system's supported page sizes in sorted order starting with the smallest supported page size. The number of array entries updated with page sizes is returned.</p> <p><b>VMINFO_PSIZE</b> Reports detailed VMM statistics for a specified page size. The <i>out</i> parameter must point to a <b>vminfo_psize</b> structure with the <i>psize</i> field set to a page size, in bytes, for which to return statistics. The <i>arg</i> parameter should be the size of the <b>vminfo_psize</b> structure.</p>
<i>arg</i>	An additional parameter that will depend upon the <i>command</i> parameter.

## Execution Environment

The **kvmgetinfo** kernel service can be called from the process environment only.

## Return Values

The following return values apply to all commands other than **VMINFO\_GETPSIZES**:

Item	Description
0	Indicates successful completion.
ENOSYS	Indicates the <i>command</i> parameter is not valid (or not yet implemented).
EINVAL	When <b>VM_PAGE_INFO</b> is the command, the <i>adr</i> field of the <b>vm_page_info</b> structure is an invalid address.

When **VMINFO\_GETPSIZES** is specified as the command, -1 is returned if the **kvmgetinfo()** kernel service is unsuccessful. Otherwise, the **kvmgetinfo()** kernel service returns a number of page sizes when the **VMINFO\_GETPSIZES** command is specified.

### Related information:

Memory Kernel Services

## kwpar\_checkpoint\_status Kernel Service

### Purpose

Provides a method for kernel services to inform the system that an event occurred within a workload partition (WPAR) that denies or subsequently reallows a checkpoint of the WPAR.

### Syntax

```
#include <sys/wparid.h>
```

```
int kwpar_checkpoint_status (kcid, cmd, varp)
```

```
cid_t kcid;
```

```
int cmd;
```

```
void * varp;
```

### Parameters

Item	Description
<i>cmd</i>	An integer command that informs the API what action to take on behalf of the caller.
<i>kcid</i>	The WPAR ID where the command operation is to take place.
<i>varp</i>	A void pointer to different elements that depends on the <i>cmd</i> parameter. <ul style="list-style-type: none"><li>• If the <i>cmd</i> parameter is set to the <b>WPAR_CHECKPOINT_TRY</b> value, the <i>varp</i> parameter is a pointer to an integer variable that contains the number of seconds that the caller is willing to wait before a blocking event is removed.</li><li>• If the <i>cmd</i> parameter is set to the <b>WPAR_CHECKPOINT_DENY</b> value, the <i>varp</i> parameter is a pointer to a null terminated character string that contains a user readable reason for posting the event.</li></ul>

### Cmd Types

The *cmd* parameter is supplied on input to the **kwpar\_checkpoint\_status** API and describes the type of action or event notification the caller is expecting. The following *cmd* types are supported:

Item	Description
<b>WPAR_CHECKPOINT_DENY</b>	The caller is experiencing an event within the WPAR identified by the <i>kcid</i> parameter that would deny a checkpoint operation. The caller must supply a pointer to a user readable character string in the <i>varp</i> parameter.
<b>WPAR_CHECKPOINT_ALLOW</b>	The caller is clearing a previous checkpoint denial operation. Deny and allow operations are cumulative and thus each denial operation must be matched with an allow operation before a checkpoint is finally reallocated.
<b>WPAR_CHECKPOINT_TRY</b>	Used by the AIX checkpoint system itself. The caller supplies the <i>varp</i> pointer to an integer that contains a "willing to wait" timeout in seconds before a checkpoint denial operation is cleared.
<b>WPAR_CHECKPOINT_CLEAR</b>	Used by the AIX checkpoint system itself. The caller completed a checkpoint after a successful <b>WPAR_TRY_CHKPNT</b> operation.
<b>WPAR_RESTART_CLEAR</b>	Used by the AIX checkpoint system itself. The caller completed a restart. The WPAR restart state is initially set when the WPAR is re-created on the arrival system.

### Description

The **kwpar\_checkpoint\_status** kernel service provides a mechanism for kernel services to inform or query the system about a checkpoint denial event. Kernel extensions that experience a temporary event which prevents a WPAR from being the target of a checkpoint operation, must use this API to deny and then to subsequently reallocate a checkpoint when the event clears. An example denial event might occur if a device open is in an unserialized interim state that cannot handle a checkpoint operation.

## Execution Environment

The `kwpar_checkpoint_status` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Success.
non-zero	Failure.

## Error Codes

The `kwpar_checkpoint_status` service fails if one or more of the following errors occur:

Item	Description
EINVAL	The caller supplied an invalid <code>cmd</code> or other parameter.
ENOENT	No WPAR with the <code>kcid</code> ID is active in the system.
EBUSY	Either of the following situations can lead to the <b>EBUSY</b> error. <ul style="list-style-type: none"><li>• WPAR is in a checkpoint or restart state. The caller is unsuccessful in a <b>WPAR_CHECKPOINT_DENY</b> operation.</li><li>• WPAR is in a state that cannot participate in a checkpoint. The caller is unsuccessful in a <b>WPAR_CHECKPOINT_TRY</b> operation.</li></ul>
ETIMEDOUT	The caller is waiting for a timeout period during a <b>WPAR_CHECKPOINT_TRY</b> operation but the timer expired.

### Related reference:

“WPAR\_CKPT\_QUERY (Checkpoint Query) Device Driver ioctl Operation” on page 584

## kwpar\_err Kernel Service

### Purpose

Logs an error message for a given Workload Partition.

### Syntax

```
int kwpar_err(kcid,cat_file_name,msg_set_no,msg_no,default_fmt_msg)
cid_t kcid;
char* cat_file_name;
unsigned int msg_set_no;
unsigned int msg_no;
char* default_fmt_msg;
```

### Description

The `kwpar_err` interface provides a mechanism to log error messages for a given WPAR from a kernel routine. Each WPAR can hold up to 1 KB of error messages. If there is enough space to log the new message, the command logs the message; otherwise, it fails. The `kwpar_err` routine is pinned and as such can be called from the interrupt handlers as well.

### Parameter

Item	Description
<i>kcid</i>	Specifies the <i>cid</i> of the WPAR.
<i>cat_file_name</i>	Specifies the catalog file name to be used for translation.
<i>msg_set_no</i>	Specifies the message set number of the error message in the catalog file.
<i>msg_no</i>	Specifies the message number of the error message.
<i>default_fmt_msg</i>	Specifies the default message string. Follows the same syntax as the <b>printf</b> subroutine <i>Format</i> parameter. Floating point is not supported.
...	Specifies the arguments to the message if any.

## Return values

Item	Description
0	Success
-1	Failure

## Error codes

Item	Description
ENOMEM	Not enough memory
EINVAL	Invalid parameter

## Example

To log an error message into WPAR with *cid* 4, enter

```
kwpar_err(4, "wparerrs.cat",1,10,"%s : command failed", "mycommand");
...
```

### Related information:

wpar\_log\_err subroutine  
wpar\_print\_err subroutine  
wparerr command

## kwpar\_getname Kernel Service

### Purpose

Returns the workload partition name associated with the requested ID.

### Syntax

```
#include<sys/wparid.h>
#include<sys/xmem.h>
```

```
int kwpar_getname(kcid, buffer, length, adspace)
cid_t kcid;
char * buffer;
size_t length;
int adspace;
```

### Description

Get the name associated with the workload partition ID (*kcid*) and write it to the output buffer. The maximum number of bytes to write is limited by the *length* parameter. The *length* parameter cannot exceed MAXCORRALNAMELEN. The service writes to either user space or kernel space, depending on the value specified for the *adspace* parameter.

### Parameters

Item	Description
<i>kcid</i>	Specifies the workload partition ID.
<i>buffer</i>	Points to the buffer where the workload partition name is stored.
<i>length</i>	Specifies the maximum number of bytes to return.
<i>adspace</i>	Indicates in which part of memory the buffer parameter is located:
	<b>SYS_ADSpace</b>
	Indicates that the <i>buffer</i> parameter is in the kernel memory.
	<b>USER_ADSpace</b>
	Indicates that the <i>buffer</i> parameter is in the application memory.

## Execution Environment

Process environment only.

## Return Values

Item	Description
0	The command completed successfully.
EINVAL	Invalid WPAR ID or specified length is greater than MAXCORRALNAMELEN.
EFAULT	Error during copyout to user space.

## kwpar\_getrootpath Kernel Service Purpose

Returns the root path of the workload partition associated with the requested ID.

## Syntax

```
#include<sys/wparid.h>
```

```
int kwpar_getrootpath(kcid, length, buffer)
cid_t kcid;
size_t * length;
char * buffer;
```

## Description

Get the root path of the workload partition associated with the *kcid* parameter and copy it to the output buffer. On entry, the value specified for the *length* parameter indicates the size of the output buffer. On return, the value specified for the *length* parameter, contains the size of the root path. If the value for the *length* parameter on entry is smaller than the actual path length, then **ENOSPC** is returned. Then, the *length* parameter is set to the actual length of the root path.

## Parameters

Item	Description
<i>kcid</i>	Specifies the workload partition ID.
<i>length</i>	Specifies the maximum number of bytes to return.
<i>buffer</i>	Points to the buffer where the workload partition root path will be stored.

## Execution Environment

Process environment only.

## Return Values

Item	Description
0	The command completed successfully.
EINVAL	Error indicating that <i>buffer</i> is NULL, <i>length</i> is NULL, or <i>*length</i> is 0.
ENOENT	Invalid WPAR ID specified for the <i>kcid</i> parameter.
ENOSPC	Insufficient space in <i>buffer</i> to copy path.

## kwpar\_isappwpar Kernel Service

### Purpose

Returns whether a workload partition is an application workload partition.

### Syntax

```
#include <sys/wparid.h>

int kwpar_isappwpar(kcid)
cid_t kcid;
```

### Description

Checks whether the workload partition associated with the *kcid* is an application workload partition.

### Parameters

Item	Description
<i>kcid</i>	Specifies the workload partition ID.

## Execution Environment

Process environment only.

## Return Values

Item	Description
1	Workload partition is an application workload partition.
0	Workload partition is not an application workload partition.
-1	Indicates that the command did not complete successfully.

## kwpar\_r2vmap\_devno Kernel Service

### Purpose

Maps a real device number to the corresponding virtual device number for a given workload partition (WPAR).

### Syntax

```
#include <sys/wparid.h>

int kwpar_r2vmap_devno ( wparid, vdevno, rdevno)
cid_t wparid;
dev_t rdevno;
dev_t * vdevno;
```

### Parameters

Item	Description
<i>wparid</i>	WPAR identifier. This parameter is required.
<i>rdevno</i>	Real device number. This parameter is required.
<i>vdevno</i>	Points to the data area that will contain the virtual device number. This parameter is passed by reference. This parameter is optional.

## Description

The **kwpar\_r2vmap\_devno** kernel service provides the ability to translate a real device number, maintained in the kernel device switch table, to the corresponding virtual device number maintained in the user space. The caller must specify an existing WPAR identifier with the *wparid* parameter and a valid real device number with the *rdevno* parameter. The **kwpar\_r2vmap\_devno** kernel service writes the corresponding virtual device number to the data area pointed to by the *vdevno* parameter (if specified). If the *vdevno* parameter is not specified, the return code indicates whether a mapping exists for the given WPAR identifier and real device number.

A mapping for the specified virtual device number must exist for the **kwpar\_v2rmap\_devno** kernel service to succeed.

## Execution Environment

The **kwpar\_r2vmap\_devno** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Success.
non-zero	Failure.

## Error Codes

The **kwpar\_r2vmap\_devno** service fails if one or more of the following errors occur:

Item	Description
EINVAL	Either the <i>wparid</i> or <i>rdevno</i> argument is invalid.
ENXIO	Unable to locate the WPAR device map associated with the given WPAR ID.
ESRCH	Unable to locate a mapping for the given real device number <i>rdevno</i> .

### Related reference:

“kwpar\_v2rmap\_devno Kernel Service” on page 326

“kwpar\_regdevno Kernel Service” on page 321

“kwpar\_unregdevno Kernel Service” on page 324

## kwpar\_r2vmap\_pid Kernel Service Purpose

Maps a real process ID to the equivalent virtual process ID assigned within a workload partition.

### Syntax

```
#include <sys/wparid.h>
```

```
pid_t kwpar_r2vmap_pid ( kcidp, rpid)
```

```
cid_t * kcidp;
```

```
pid_t rpid;
```

## Parameters

Item	Description
<i>kcidp</i>	A pointer to a memory location where the workload partition (WPAR) ID associated with the <i>rpid</i> parameter is returned.
<i>rpid</i>	The real process ID on which to translate a real process ID to a virtual process ID.

## Description

The **kwpar\_r2vmap\_pid** kernel service provides a mapping from a real process ID to a virtual process ID assigned within the workload partition. In most instances, the real and virtual process IDs are the same except in cases where the Workload Partition Mobility is in effect or for certain system services such as the **init** command which always have different real and virtual process IDs.

Usually kernel services dealing with process IDs only accept real process IDs. However, in some instances it might be necessary for kernel extensions, which communicate with other WPAR services or with processes within the WPAR, to know and communicate with virtual process IDs.

## Execution Environment

The **kwpar\_r2vmap\_pid** kernel service can be called from the process environment only.

## Return Values

If the **kwpar\_r2vmap\_pid** kernel service succeeds, it returns the virtual *pid\_t* value associated with the *rpid* value provided on input. If the kernel service fails or if there is no virtual process ID associated with the *rpid* value, the *rpid* value is returned.

### Related reference:

“kwpar\_v2rmap\_pid Kernel Service” on page 327

## kwpar\_r2vmap\_tid Kernel Service

### Purpose

Maps a real thread ID to the equivalent virtual thread ID assigned within a workload partition.

### Syntax

```
#include <sys/wparid.h>
```

```
tid_t kwpar_r2vmap_tid ( kcidp, rtid)
```

```
cid_t * kcidp;
```

```
tid_t rtid;
```

## Parameters

Item	Description
<i>kcidp</i>	A pointer to a memory location where the WPAR ID associated with the <i>rtid</i> parameter is returned.
<i>rtid</i>	The real thread ID on which to translate a real process ID to a virtual process ID.

## Description

The **kwpar\_r2vmap\_tid** kernel service provides a mapping from a real thread ID to a virtual thread ID assigned within the workload partition. In most instances, the real and virtual thread IDs are the same except in cases where the Workload Partition Mobility is in effect.



Normally kernel services dealing with thread IDs accept only real thread IDs. However, in some instances it might be necessary for kernel extensions, which communicate with other WPAR services or with processes within the WPAR, to know and communicate with virtual thread IDs.

## Execution Environment

The `kwpar_r2vmap_tid` kernel service can be called from the process environment only.

## Return Values

If the `kwpar_r2vmap_tid` kernel service succeeds, it returns the virtual `tid_t` value associated with the `rtid` value provided on input. If the kernel service fails or if there is no virtual process ID associated with the `rtid` value, the `rtid` value is returned.

### Related reference:

“`kwpar_v2rmap_tid` Kernel Service” on page 328

## kwpar\_regdevno Kernel Service

### Purpose

Registers a virtual device number for a given workload partition (WPAR) by mapping it to a real device number in the device switch table.

### Syntax

```
#include <sys/wparid.h>
```

```
int kwpar_regdevno ( wparid, vdevno, rdevno)
cid_t wparid;
dev_t vdevno;
dev_t * rdevno;
```

### Parameters

Item	Description
<code>wparid</code>	WPAR ID. This parameter is required.
<code>vdevno</code>	Virtual device number. This parameter is required.
<code>rdevno</code>	Points to the data area that will contains the real device number. This parameter is passed by reference. This parameter is required.

### Description

The `kwpar_regdevno` kernel service provides the ability to register a virtual device number for a given WPAR by mapping it to a real device number in the device switch table. The `kwpar_regdevno` kernel service performs the following steps:

1. Locates a free slot in the kernel device switch table and reserves it for the WPAR specified by the `wparid` parameter.
2. Creates a mapping between the virtual device number, which is specified by the `vdevno` parameter, to the real device number reserved in the previous step.
3. The newly reserved real device number is passed back to the caller through the `rdevno` parameter.

## Execution Environment

The `kwpar_regdevno` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Success.
non-zero	Failure.

## Error Codes

The `kwpar_regdevno` kernel service fails if one or more of the following errors occur:

Item	Description
EINVAL	Either the <code>wparid</code> or <code>vdevno</code> argument is not valid.
ENXIO	Unable to locate the WPAR device map associated with the given WPAR ID.
ENOTEMPTY	The virtual device number <code>vdevno</code> is already mapped.

### Related reference:

“`kwpar_r2vmap_devno` Kernel Service” on page 318

“`kwpar_v2rmap_devno` Kernel Service” on page 326

“`kwpar_unregdevno` Kernel Service” on page 324

## `kwpar_reghook` Kernel Service Purpose

Registers a function callback with workload partition (WPAR) kernel services. Callback functions are subsequently performed when specific WPAR conditions occur.

### Syntax

```
#include <sys/wparid.h>
```

```
regkey_t kwpar_reghook ( hooktype, hookp)
int hooktype;
void * hookp;
```

### Parameters

Item	Description
<code>hooktype</code>	Identifies the form of the <code>hookp</code> pointer.
<code>hookp</code>	A pointer to a memory location that might contain function pointers or other structure elements that are interpreted depending on the supplied <code>hooktype</code> value.

### Hook Types

The `hooktype` parameter is supplied on input to the `kwpar_reghook` return and describes the form of the second parameter. The supported hook types are as follows:

Item	Description
WPAR_NOTIFY_HOOK	Identifies the form of the <code>hookp</code> parameter as being of type <code>wpar_config_hook_t</code> .

The `wpar_config_hook_t` structure contains the following fields:

Item	Description
<code>uint current_hiwater</code>	On output from the <code>kwpar_reghook</code> service, this field contains the current upper number of WPARs that became active on this boot instance of the AIX operating system. WPAR IDs are allocated in numeric order. Kernel subsystems that want to size internal components according to the number of active WPARs must register a <code>WPAR_NOTIFY_HOOK</code> hook type and examine the <code>current_hiwater</code> value for existing WPARs during registration. Future WPAR activation after hook registration calls the specified <code>configp</code> function within the <code>wpar_config_hook_t</code> element. See the <code>WPARSTART</code> flags later in this section for a further description of the WPAR activation.
<code>wpar_config_func_t configp</code>	On input, this field contains a pointer to a callback routine that is started by the WPAR kernel services during the activation and the deactivation of workload partitions within the AIX kernel.

The syntax for the `wpar_config_func_t` is as follows:

```
#include <sys/wpar.h>
```

```
typedef int * wpar_config_func_t ( flags, cid, corralp, unused)
int flags;
cid_t cid;
struct corral * corralp;
void * unused;
```

The parameters are as follows:

Item	Description
<code>flags</code>	Information regarding the type of condition that is occurring within the workload partition.
<code>cid</code>	The ID for the workload partition experiencing the condition.
<code>corralp</code>	A pointer to a kernel copy of the <code>corral</code> structure that might be supplied from the user space at the start of the condition processing.
<code>unused</code>	Currently unused and must be set to NULL. It might be expanded to contain more information in later revisions of this API.

The `flags` parameter can have the following potential values:

Item	Description
<code>WPARSTART</code>	Signifies that the WPAR is undergoing activation. The callout to registered routines occurs before any other kernel subsystem processing occurs. Kernel components registering and desiring to see the WPAR activation are informed that a new WPAR with the <code>cid</code> parameter set is going to enter the AIX kernel system.
<code>WPARSTOP</code>	Signifies that the WPAR underwent deactivation. The callout to registered routines occurs after all other kernel subsystem processing occurs. Kernel components registering and desiring to see the WPAR deactivation are informed that an existing WPAR with the <code>cid</code> parameter set left the AIX kernel system.

## Description

The `kwpar_reghook` kernel service provides a mechanism for other kernel services to register callbacks and retrieve information when certain workload partition conditions occur.

## Execution Environment

The `kwpar_reghook` kernel service can be called from the process environment only.

## Return Values

If the `kwpar_reghook` kernel service is successful, it returns a registration key that can subsequently be used with the `kwpar_unreghook` kernel service. If the kernel service fails, it returns a numeric value equivalent to the `BADREGKEY` definition found in the `wparid.h` file.

## Error Codes

The `kwpar_reghook` kernel service fails if no space remains to record additional registration hook.

**Related reference:**

“`kwpar_unreghook` Kernel Service” on page 325

## `kwpar_unregdevno` Kernel Service

### Purpose

Unregisters the mapping associated with a real device number for a given workload partition (WPAR).

### Syntax

```
#include <sys/wparid.h>
```

```
int kwpar_unregdevno ( wparid, rdevno)  
cid_t wparid;  
dev_t rdevno;
```

### Parameters

Item	Description
<i>wparid</i>	WPAR identifier. This parameter is required.
<i>rdevno</i>	Real device number. This parameter is required.

### Description

The `kwpar_unregdevno` kernel service provides the ability to unregister the mapping associated with a real device number for a given WPAR. The `kwpar_unregdevno` kernel service will perform the following steps:

1. Deletes the virtual-to-real mapping associated with the real device number specified by the *rdevno* parameter for the WPAR specified by the *wparid* parameter.
2. Releases the reserve associated with the real device number specified by the *rdevno* parameter.

### Execution Environment

The `kwpar_unregdevno` kernel service can be called from the process environment only.

### Return Values

Item	Description
0	Success.
non-zero	Failure.

## Error Codes

The `kwpar_unregdevno` kernel service fails if one or more of the following errors occur:

Item	Description
EINVAL	Either the <i>wparid</i> or <i>rdevno</i> argument is not valid.
ENXIO	Unable to locate the WPAR device map associated with the given WPAR ID.
ESRCH	Unable to locate the mapping for the given real device number <i>rdevno</i> .

#### Related reference:

“*kwpar\_r2vmap\_devno* Kernel Service” on page 318

“*kwpar\_v2rmap\_devno* Kernel Service” on page 326

“*kwpar\_regdevno* Kernel Service” on page 321

## kwpar\_unreghook Kernel Service

### Purpose

Removes a previously registered workload partition (WPAR) callback hook.

### Syntax

```
#include <sys/wparid.h>
```

```
int kwpar_unreghook ( key)
regkey_t key;
```

### Parameters

Item	Description
<i>key</i>	The registration key of the hook that the caller wants to un-register. This key is equivalent to the key returned from a hook registration with the <b>kwpar_reghook</b> kernel service.

### Description

The **kwpar\_unreghook** kernel service informs workload partitions that the caller no longer wants to receive callouts for WPAR conditions.

### Execution Environment

The **kwpar\_unreghook** kernel service can be called from the process environment only.

### Return Values

Item	Description
0	Success.
non-zero	Failure.

### Error Codes

The **kwpar\_unreghook** service fails if one or more of the following errors occur:

Item	Description
EINVAL	Not a valid registration key.
EPERM	Not allowed to un-register this key.

#### Related reference:

“kwpar\_reghook Kernel Service” on page 322

## kwpar\_v2rmap\_devno Kernel Service

### Purpose

Maps a virtual device number to the corresponding real device number in the device switch table for a given workload partition (WPAR).

### Syntax

```
#include <sys/wparid.h>
```

```
int kwpar_v2rmap_devno ( wparid, vdevno, rdevno)
cid_t wparid;
dev_t vdevno;
dev_t * rdevno;
```

### Parameters

Item	Description
<i>wparid</i>	WPAR identifier. This parameter is required.
<i>vdevno</i>	Virtual device number. This parameter is required.
<i>rdevno</i>	Points to the data area that will contain the real device number. This parameter is passed by reference. This parameter is optional.

### Description

The **kwpar\_v2rmap\_devno** kernel service provides the ability to translate a virtual device number maintained in user space to the corresponding real device number maintained in the kernel device switch table. The caller must specify an existing WPAR identifier with the *wparid* parameter and a valid virtual device number with the *vdevno* parameter. The **kwpar\_v2rmap\_devno** kernel service will write the corresponding real device number to the data area pointed to by the *rdevno* parameter if it is specified. If the *rdevno* parameter is not specified, the return code will indicate whether a mapping exists for the given WPAR identifier and virtual device number.

A mapping for the specified virtual device number must exist for the **kwpar\_v2rmap\_devno** kernel service to succeed.

### Execution Environment

The **kwpar\_v2rmap\_devno** kernel service can be called from the process environment only.

### Return Values

Item	Description
0	Success.
non-zero	Failure.

## Error Codes

The `kwpar_v2rmap_devno` service fails if one or more of the following errors occur:

Item	Description
EINVAL	Either the <code>wparid</code> or <code>vdevno</code> argument is not valid.
ENXIO	Unable to locate the WPAR device map associated with the given WPAR id.
ENODEV	Unable to locate the mapping for the given virtual device number.

### Related reference:

“`kwpar_r2vmap_devno` Kernel Service” on page 318

“`kwpar_regdevno` Kernel Service” on page 321

“`kwpar_unregdevno` Kernel Service” on page 324

## `kwpar_v2rmap_pid` Kernel Service Purpose

Maps a virtual process ID associated with a process within a workload partition to the equivalent real process ID.

### Syntax

```
#include <sys/wparid.h>
```

```
pid_t kwpar_v2rmap_pid ( kcid, vpid)
cid_t kcid;
pid_t vpid;
```

### Parameters

Item	Description
<code>kcid</code>	The workload partition (WPAR) ID associated with the <code>vpid</code> parameter. Equivalent virtual process IDs can be in use across different processes in different WPARs. Thus the caller must provide the WPAR ID for which a virtual to real mapping is to occur.
<code>vpid</code>	The virtual process ID on which to perform a virtual to real mapping.

### Description

The `kwpar_v2rmap_pid` kernel service provides a mapping from a virtual process ID associated with a process in a workload partition to the equivalent real process ID. In most instances, both the real and virtual process IDs are the same, except in cases where the Workload Partition Mobility is in effect.

Normally, kernel services dealing with process IDs accept only real thread IDs. In some instances where a kernel extension is communicating with other WPAR services or with processes within the WPAR, a mapping from virtual to real process IDs might be needed.

### Execution Environment

The `kwpar_v2rmap_pid` kernel service can be called from the process environment only.

## Return Values

If the `kwpar_v2rmap_pid` kernel service succeeds, it returns the real `pid_t` value associated with the `vpid` value provided on input. If the kernel service fails, or if there is no real thread ID associated with the `vpid` value, then the `vpid` value is returned.

### Related reference:

“`kwpar_r2vmap_pid` Kernel Service” on page 319

## `kwpar_v2rmap_tid` Kernel Service

### Purpose

Maps a virtual thread ID associated with a thread within a workload partition to the equivalent real thread ID.

### Syntax

```
#include <sys/wparid.h>
```

```
tid_t kwpar_v2rmap_tid ( kcid, vtid)
```

```
cid_t kcid;
```

```
tid_t vtid;
```

### Parameters

Item	Description
<code>kcid</code>	The workload partition (WPAR) ID associated with the <code>vtid</code> parameter. Equivalent virtual thread IDs can be in use across different threads in different WPARs. Thus the caller must provide the WPAR ID for which a virtual to real mapping is to occur.
<code>vtid</code>	The virtual thread ID on which to perform a virtual to real mapping.

### Description

The `kwpar_v2rmap_tid` kernel service provides a mapping from a virtual thread ID associated with a thread in a workload partition to the equivalent real thread ID. In most instances, both the real and virtual thread IDs are the same, except in cases where the Workload Partition Mobility is in effect. Normally, kernel services dealing with thread IDs accept only real thread IDs. In some instances where a kernel extension is communicating with other WPAR services or with processes within the WPAR, a mapping from virtual to real thread IDs might be needed.

### Execution Environment

The `kwpar_v2rmap_tid` kernel service can be called from the process environment only.

### Return Values

If the `kwpar_v2rmap_tid` kernel service succeeds, it returns the real `tid_t` value associated with the `vtid` value provided on input. If the kernel service fails, or if there is no real thread ID associated with the `vtid` value then the `vtid` value is returned.

### Related reference:

“`kwpar_r2vmap_tid` Kernel Service” on page 320

## I

The following kernel services begin with the with the letter I.



## **ldata\_alloc Kernel Service**

### **Purpose**

Allocates a pinned storage element from an **ldata** pool.

### **Syntax**

```
#include <sys/ldata.h>
```

```
void * ldata_alloc (ldatap)
```

```
ldata_t ldatap;
```

### **Description**

The **ldata\_alloc** kernel service allocates a pinned storage element from a **ldata** pool and returns the address of the element. The **ldata\_alloc** kernel service makes a pinned storage element from the **ldata** pool available for use by the caller. The sub-pool from which the element is allocated corresponds to the SRAD on which the call was made. If there are no free pinned elements, a new element cannot be allocated and a NULL value is returned.

After it is allocated, the pinned storage element can be freed to the **ldata** pool through the **ldata\_free** kernel service.

### **Parameters**

<b>Item</b>	<b>Description</b>
<i>ldatap</i>	Specifies the handle of the <b>ldata</b> pool.

### **Execution Environment**

The **ldata\_alloc** kernel service can be called from the process or interrupt environment.

### **Return Values**

Returns a pointer to a pinned storage element allocated from an **ldata** pool or NULL if no element could be allocated.

### **Implementation Specifics**

The **ldata\_alloc** kernel service is part of the Base Operating System (BOS) Runtime.

#### **Related reference:**

“ldata\_create Kernel Service”

“ldata\_grow Kernel Service” on page 332

“ldata\_free Kernel Service” on page 332

## **ldata\_create Kernel Service**

### **Purpose**

Creates a SRAD-aware pinned storage element pool (**ldata** pool) and returns its handle.

### **Syntax**

```
#include <sys/ldata.h>
```

```
int ldata_create (size, initcount, maxcount, kkey, ldatap)
```

```

size_t size;
long initcount;
long maxcount;
kkey_t kkey;
ldata_t * ldatap;

```

## Description

The **ldata\_create** kernel service creates a SRAD-aware pool (**ldata** pool) of pinned storage elements, each of the specified size, and returns a handle to the newly-allocated pool. An **ldata** pool consists of a number of sub-pools (one per SRAD). Each sub-pool is physically backed with memory local to its corresponding SRAD. The size of each sub-pool is equal to the value of the *maxcount* parameter multiplied by the value of the *size* parameter. The parameter (*initcount*) specifies the number of pinned storage elements in each sub-pool that should be pre-allocated.

The **ldata** pool can be created with a kernel storage protection key by specifying one through the *kkey* parameter. For compatibility with previous releases, a *kkey* parameter of zero requests no protection. When a protection key is specified, the caller must hold this key when calling any **ldata** service, including the **ldata\_create** kernel service.

After an **ldata** pool is created, its handle can be used to allocate pinned storage elements from the pool through the **ldata\_alloc** kernel service and free these elements to the pool through the **ldata\_free** kernel services. Elements are allocated and freed to the sub-pool corresponding to the SRAD on which **ldata\_alloc** and **ldata\_free** are called. If a sub-pool is exhausted of its pinned storage elements, it can be grown by calling the **ldata\_grow** kernel service up to *maxcount*.

An **ldata** pool created through the **ldata\_create** service can be destroyed by the **ldata\_destroy** kernel service.

## Parameters

Item	Description
<i>size</i>	Specifies the size, in bytes, of each pinned storage element of the <b>ldata</b> pool.
<i>initcount</i>	Specifies the initial count of pinned storage elements, to be contained within the <b>ldata</b> pool. Must be a positive integer.
<i>maxcount</i>	Specifies the maximum count of pinned storage elements that can be contained with the <b>ldata</b> pool. The value of <i>maxcount</i> must be positive and greater than or equal to the value of <i>initcount</i> .
<i>kkey</i>	Specifies the kernel storage protection key to be applied to the newly created <b>ldata</b> pool. The value must be a valid kernel key number, or zero to indicate that storage protection is not requested.
<i>ldatap</i>	Specifies an address to be set on successful completion with the handle for the newly created <b>ldata</b> pool.

## Execution Environment

The **ldata\_create** kernel service can be called only from the process environment.

## Return Values

Item	Description
0	Completed successfully. The handle for <b>ldata</b> storage is returned in <i>ldatap</i> .
EINVAL	Invalid input parameters given. Invalid <i>initcount</i> , <i>maxcount</i> or <i>kkey</i> . The <i>ldatap</i> parameter is undefined.
ENOMEM	Error encountered. Insufficient memory to satisfy request. The <i>ldatap</i> parameter is undefined.

## Implementation Specifics

The **ldata\_create** kernel service is part of the Base Operating System (BOS) Runtime.

### Related reference:

“ldata\_destroy Kernel Service”

“ldata\_grow Kernel Service” on page 332

“ldata\_alloc Kernel Service” on page 329

## ldata\_destroy Kernel Service

### Purpose

Destroys an **ldata** pool created by the **ldata\_create** kernel service.

### Syntax

```
#include <sys/ldata.h>
```

```
void ldata_destroy (ldatap)
```

```
ldata_t ldatap;
```

### Description

The **ldata\_destroy** kernel service destroys an **ldata** pool previously created by an **ldata\_create** call. This routine assumes that all elements allocated from the pool have been freed back to the pool and there are no longer any active elements in the pool.

The **ldata\_destroy** call unpins and frees all of the storage associated with the handle.

### Parameters

Item	Description
<i>ldatap</i>	Specifies the handle of the <b>ldata</b> pool to be destroyed.

### Execution Environment

The **ldata\_destroy** kernel service can be called from the process environment only.

### Return Values

None.

## Implementation Specifics

The **ldata\_destroy** kernel service is part of the Base Operating System (BOS) Runtime.

### Related reference:

“ldata\_create Kernel Service” on page 329

“ldata\_alloc Kernel Service” on page 329

“ldata\_free Kernel Service”

## **ldata\_free Kernel Service**

### **Purpose**

Frees a storage element that is pinned to an **ldata** pool.

### **Syntax**

```
#include <sys/ldata.h>
```

```
void ldata_free (ldatap, elementp)
```

```
ldata_t ldatap;  
void * elementp;
```

### **Description**

The **ldata\_free** kernel service frees a pinned storage element that was previously allocated to an **ldata** pool. The pinned storage element is identified through the *elementp* parameter. The element identified by *elementp* is freed to the sub-pool corresponding to the SRAD that allocated the element.

### **Parameters**

Item	Description
<i>ldatap</i>	Specifies the handle of the <b>ldata</b> pool.
<i>elementp</i>	Specifies the address of the pinned storage element to be freed.

### **Execution Environment**

The **ldata\_free** kernel service can be called from the process or interrupt environment.

### **Return Values**

None.

### **Implementation Specifics**

The **ldata\_free** kernel service is part of Base Operating System (BOS) Runtime.

#### **Related reference:**

“ldata\_alloc Kernel Service” on page 329

## **ldata\_grow Kernel Service**

### **Purpose**

Expands the count of available pinned storage elements contained within an **ldata** pool.

### **Syntax**

```
#include <sys/ldata.h>
```

```
int ldata_grow (ldatap, count)
```

```
ldata_t ldatap;  
long count;
```

## Description

The `ldata_grow` kernel service increases the number of pinned storage elements contained within a per-SRAD sub-pool associated with the `ldata` handle `ldatap`, by `count`. If the `ldata_alloc` call fails because there are no more free pinned storage elements in a sub-pool, use the `ldata_grow` kernel service. The `ldata_grow` kernel service pins additional count elements from the sub-pool and makes them available for the `ldata_alloc` call. All of the sub-pools associated with the handle are grown. If count elements are not available or there is not enough pinned memory available, the `ldata_grow` kernel service fails.

## Parameters

Item	Description
<code>ldatap</code>	Specifies the handle of the <code>ldata</code> pool.
<code>count</code>	Specifies the additional number of storage elements to be pinned in the sub-pool. The <code>count</code> value should be greater than 0 and should not increase the sub-pool size beyond the value of <code>maxcount</code> specified with the <code>ldata_create</code> call.

## Execution Environment

The `ldata_grow` kernel service can be called only from the process environment.

## Return Values

Item	Description
0	Success.
-1	Error encountered. Illegal parameters or insufficient resources.

## Implementation Specifics

The `ldata_grow` kernel service is part of the Base Operating System (BOS) Runtime.

Related reference:

“`ldata_create` Kernel Service” on page 329

## `ldmp_bufest`, `ldmp_timeleft`, `ldmp_xmalloc`, `ldmp_xmfree`, and `ldmp_errstr` Kernel Services

### Purpose

Obtains information about the current live dump.

### Syntax

```
#include <sys/livedump.h>
```

```
kerrno_t ldmp_bufest (id, cb, len)
dumpid_t id;
ras_block_t cb;
size_t *len;
```

```
kerrno_t ldmp_timeleft (id, timeleft)
dumpid_t id;
long *timeleft;
```

```
kerrno_t ldmp_xmalloc (id, size, align, p)
dumpid_t id;
```

```
size_t size;
uint align;
void **p;
```

```
kernno_t ldmp_xmfree (id, p)
dumpid_t id;
void *p;
```

```
kernno_t ldmp_errstr (id, cb, str)
dumpid_t id;
ras_block_t cb;
char *str;
```

## Parameters

Item	Description
<i>align</i>	Specifies the log base 2 of the desired alignment. The maximum allowed alignment is 12, 4096 byte alignment.
<i>cb</i>	Specifies the <code>ras_block_t</code> for the component.
<i>id</i>	Specifies the ID of the dump.
<i>len</i>	Specifies the estimate of data in bytes that can still be buffered by the specified component in this pass.
<i>p</i>	Specifies the memory block to be allocated or freed.
<i>size</i>	Specifies the memory size to be allocated.
<i>str</i>	Specifies the error message.
<i>timeleft</i>	Specifies the time, in nanoseconds, remaining for this pass. This value only has meaning for a serialized dump. It can be negative.

## Description

The `ldmp_bufest` kernel service estimates the number of bytes of dump buffer storage available to this component.

The `ldmp_timeleft` kernel service estimates the time, in nanoseconds, remaining in this pass.

The `ldmp_xmalloc` kernel service allocates storage from the live dump heap.

The `ldmp_xmfree` kernel service frees live dump heap storage.

The `ldmp_errstr` kernel service records an error to be part of the live dump status reporting. The string is contained in the live dump and reported in the error log entry if there is sufficient space.

**Important:** An error log entry has a maximum length of 2048 bytes. The error string is limited to 128 bytes, including the trailing NULL, and is truncated if too long. The component's path name is also logged.

**Tip:** The `ldmp_errstr` kernel service can be called multiple times to report multiple errors. Components are encouraged to limit the size of error strings due to limited space in the error log entry.

## Return Values

Item	Description
0	Indicates a successful completion.
EINVAL_RAS_xxx_BADARGS	Indicates that the arguments for the service are not valid.
EFAULT_RAS_xxx_BADARGS	Indicates that an address argument is not a valid address.
ENOMEM_RAS_LDMP_XMALLOC	Indicates that there is insufficient space in the live dump heap to satisfy this request.

**Related reference:**

“livedump Kernel Service” on page 337

## ldmp\_freeparms Kernel Service Purpose

Frees any data allocated by the live dump associated with an unused **ldmp\_parms\_t** data item.

### Syntax

```
#include <sys/livedump.h>
```

```
kernno_t ldmp_freeparms (parms)
ldmp_parms_t *parms;
```

### Parameters

Item	Description
<i>parms</i>	Points to an item of <b>ldmp_parms_t</b> type.

### Description

The **ldmp\_freeparms** kernel service is used in the event that you have partially set up the **ldmp\_parms\_t** data item, but do not want to take a dump. You can use the **ldmp\_freeparms** kernel service to clean up any data allocated by the live dump subsystem. However, you can always call the **ldmp\_freeparms** kernel service after the **livedump** kernel service, and the **ldmp\_freeparms** kernel service returns normally if there is nothing to free.

### Execution Environment

The **ldmp\_freeparms** kernel service can be called from either the process or interrupt environment.

### Return Values

Item	Description
0	Indicates a successful completion.
EINVAL_RAS_LDMP_FREEPARMS	Indicates that the area is not a valid <b>ldmp_parms_t</b> data area.
EFAULT_RAS_LDMP_FREEPARMS	Indicates that a memory fault results.

**Related reference:**

“ldmp\_setupparms Kernel Service”

“livedump Kernel Service” on page 337

## ldmp\_setupparms Kernel Service Purpose

Sets up the **ldmp\_parms\_t** parameter for the **livedump** kernel service.

## Syntax

```
#include <sys/livedump.h>
```

```
kerrno_t ldmp_setupparms (parms)  
ldmp_parms_t *parms;
```

## Parameters

Item	Description
<i>parms</i>	Points to an item of <b>ldmp_parms_t</b> type.

## Description

The **ldmp\_setupparms** kernel service simplifies the process of setting up a live dump by setting up the **ldmp\_parms\_t** parameter. It does not allocate any storage.

The **ldmp\_setupparms** kernel service performs the following setup for the **ldmp\_parms\_t** parameter:

Item	Description
Field	Value
<i>ldp_eyec</i>	eyecatcher for <b>ldmp_parms</b>
<i>ldp_vers</i>	current version
<i>ldp_flags</i>	0
<i>ldp_prio</i>	LDPP_CRITICAL
<i>ldp_recov</i>	NULL
<i>ldp_func</i>	NULL
<i>ldp_namepref</i>	NULL
<i>ldp_errcode</i>	0
<i>ldp_symptom</i>	NULL
<i>ldp_title</i>	NULL
<i>ldp_rsvd1</i>	NULL

## Execution Environment

The **ldmp\_setupparms** kernel service can be called from either the process or interrupt environment.

## Return Values

Item	Description
0	Indicates a successful completion.
EFAULT_RAS_LDMP_SETUPPARMS	Indicates that the address is not valid.

### Related reference:

“livedump Kernel Service” on page 337

## limit\_sigs or sigsetmask Kernel Service Purpose

Changes the signal mask for the calling kernel thread.

## Syntax

```
#include <sys/encap.h>
```

```
void limit_sigs (  
    siglist,
```



```

    old_mask)
sigset_t *siglist;
sigset_t *old_mask;

void sigsetmask ( old_mask)
sigset_t *old_mask;

```

## Parameters

Item	Description
<i>siglist</i>	Specifies the signal set to deliver.
<i>old_mask</i>	Points to the old signal set.

## Description

The **limit\_sigs** kernel service changes the signal mask for the calling kernel thread such that only the signals specified by the *siglist* parameter will be delivered, unless they are currently being blocked or ignored.

The old signal mask is returned via the *old\_mask* parameter. If the *siglist* parameter is **NULL**, the signal mask is not changed; it can be used for getting the current signal mask.

The **sigsetmask** kernel service should be used to restore the set of blocked signals for the calling thread. The typical usage of these services is the following:

```

sigset_t allowed = limited set of signals
sigset_t old;

/* limits the set of delivered signals */
limit_sigs (&allowed, &old);

    /* do something with a limited set of delivered signals */

/* restore the original set */
sigsetmask (&old);

```

## Execution Environment

The **limit\_sigs** and **sigsetmask** kernel services can be called from the process environment only.

## Return Values

The **limit\_sigs** and **sigsetmask** kernel services have no return values.

### Related reference:

“kthread\_kill Kernel Service” on page 310

### Related information:

Process and Exception Management Kernel Services

## livedump Kernel Service

### Purpose

Starts a live dump.

### Syntax

```
#include <sys/livedump.h>
```

```
kernno_t livedump (parms)
ldmp_parms_t *parms;
```

## Parameters

Item	Description
<i>parms</i>	Points to an item of <b>ldmp_parms_t</b> type.

## Description

The **livedump** kernel service initiates a live dump. It can be called from either the kernel or a kernel extension. Storage associated with the dump is not entirely freed until the dump has been written to disk, or the **livedump** kernel service returns an error indicating the dump was not taken.

## Execution Environment

The **livedump** kernel service can be called from either the process or interrupt environment. Only a serialized, synchronous dump can be started from the interrupt level, and the dump is limited to one pass.

## Return Values

Item	Description
0	Indicates a successful completion.
EINVAL_RAS_LIVEDUMP_PARM	Indicates that one or more parameters are not valid.
EFAULT_RAS_LIVEDUMP_PARM	Indicates that a memory fault occurs.
EINVAL_RAS_LIVEDUMP_COMP	Indicates one or more components are not valid.
EINVAL_RAS_LIVEDUMP_NOCOMPS	Indicates that no valid components were given.

## Related reference:

“dmp\_compspec and dmp\_compext Kernel Services” on page 93

“dmp\_eaddr, dmp\_context, dmp\_tid, dmp\_pid, dmp\_errbuf, dmp\_mtrc, dmp\_systrace, and dmp\_ct Kernel Services” on page 100

## lock\_alloc Kernel Service

### Purpose

Allocates system memory for a simple or complex lock.

### Syntax

```
#include <sys/lock_def.h>
#include <sys/lock_alloc.h>
```

```
void lock_alloc ( lock_addr, flags, class, occurrence)
void *lock_addr;
int flags;
short class;
short occurrence;
```

### Parameters

<b>Item</b>	<b>Description</b>
<i>lock_addr</i>	Specifies a valid simple or complex lock address.
<i>flags</i>	Specifies whether the memory allocated is to be pinned or pageable. Set this parameter as follows: <ul style="list-style-type: none"> <li><b>LOCK_ALLOC_PIN</b> Allocate pinned memory; use if it is not permissible to take a page fault while calling a locking kernel service for this lock.</li> <li><b>LOCK_ALLOC_PAGED</b> Allocate pageable memory; use if it is permissible to take a page fault while calling a locking kernel service for this lock.</li> </ul>
<b>Item</b>	<b>Description</b>
<i>class</i>	Specifies the family which the lock belongs to.
<i>occurrence</i>	Identifies the instance of the lock within the family. If only one instance of the lock is defined, this parameter should be set to -1.

## Description

The **lock\_alloc** kernel service allocates system memory for a simple or complex lock. The **lock\_alloc** kernel service must be called for each simple or complex before the lock is initialized and used. The memory allocated is for internal lock instrumentation use, and is not returned to the caller; no memory is allocated if instrumentation is not used.

## Execution Environment

The **lock\_alloc** kernel service can be called from the process environment only.

## Return Values

The **lock\_alloc** kernel service has no return values.

### Related reference:

“lock\_free Kernel Service” on page 341

“lock\_init Kernel Service” on page 342

### Related information:

Understanding Locking

## lock\_clear\_recursive Kernel Service

### Purpose

Prevents a complex lock from being acquired recursively.

### Syntax

```
#include <sys/lock_def.h>
```

```
void lock_clear_recursive ( lock_addr)
complex_lock_t lock_addr;
```

### Parameter

Item	Description
<i>lock_addr</i>	Specifies the address of the lock word which is no longer to be acquired recursively.

## Description

The **lock\_clear\_recursive** kernel service prevents the specified complex lock from being acquired recursively. The lock must have been made recursive with the **lock\_set\_recursive** kernel service. The calling thread must hold the specified complex lock in write-exclusive mode.

## Execution Environment

The **lock\_clear\_recursive** kernel service can be called from the process environment only.

## Return Values

The **lock\_clear\_recursive** kernel service has no return values.

### Related reference:

“lock\_init Kernel Service” on page 342

“lock\_done Kernel Service”

### Related information:

Locking Kernel Services

## lock\_done Kernel Service

### Purpose

Unlocks a complex lock.

### Syntax

```
#include <sys/lock_def.h>
```

```
void lock_done ( lock_addr)
complex_lock_t lock_addr;
```

### Parameter

Item	Description
<i>lock_addr</i>	Specifies the address of the lock word to unlock.

## Description

The **lock\_done** kernel services unlocks a complex lock. The calling kernel thread must hold the lock either in shared-read mode or exclusive-write mode. If one or more kernel threads are waiting to acquire the lock in exclusive-write mode, one of these kernel threads (the one with the highest priority) is made runnable and may compete for the lock. Otherwise, any kernel threads which are waiting to acquire the lock in shared-read mode are made runnable. If there was at least one kernel thread waiting for the lock, the priority of the calling kernel thread is recomputed.

If the lock is held recursively, it is not actually released until the **lock\_done** kernel service has been called once for each time that the lock was locked.

## Execution Environment

The **lock\_done** kernel service can be called from the process environment only.

## Return Values

The `lock_done` kernel service has no return values.

### Related reference:

“`lock_alloc` Kernel Service” on page 338

“`lock_free` Kernel Service”

“`lock_init` Kernel Service” on page 342

### Related information:

Understanding Locking

Locking Kernel Services

## lock\_free Kernel Service

### Purpose

Frees the memory of a simple or complex lock.

### Syntax

```
#include <sys/lock_def.h>
#include <sys/lock_alloc.h>
```

```
void lock_free ( lock_addr)
void *lock_addr;
```

### Parameter

Item	Description
<i>lock_addr</i>	Specifies the address of the lock word whose memory is to be freed.

### Description

The `lock_free` kernel service frees the memory of a simple or complex lock. The memory freed is the internal operating system memory which was allocated with the `lock_alloc` kernel service.

**Note:** It is only necessary to call the `lock_free` kernel service when the memory that the corresponding lock was protecting is released. For example, if you allocate memory for an i-node which is to be protected by a lock, you must allocate and initialize the lock before using it. The memory may be used with several i-nodes, each taken from, and returned to, the free i-node pool; the `lock_init` kernel service must be called each time this is done. The `lock_free` kernel service must be called when the memory allocated for the inode is finally freed.

### Execution Environment

The `lock_free` kernel service can be called from the process environment only.

### Return Values

The `lock_free` kernel service has no return values.

### Related reference:

“`lock_alloc` Kernel Service” on page 338

### Related information:

Understanding Locking

Locking Kernel Services

## lock\_init Kernel Service

### Purpose

Initializes a complex lock.

### Syntax

```
#include <sys/lock_def.h>
```

```
void lock_init ( lock_addr, can_sleep)  
complex_lock_t lock_addr;  
boolean_t can_sleep;
```

### Parameters

Item	Description
<i>lock_addr</i>	Specifies the address of the lock word.
<i>can_sleep</i>	This parameter is ignored.

### Description

The **lock\_init** kernel service initializes the specified complex lock. This kernel service must be called for each complex lock before the lock is used. The complex lock must previously have been allocated with the **lock\_alloc** kernel service. The *can\_sleep* parameter is included for compatibility with OSF/1 1.1, but is ignored. Using a value of **TRUE** for this parameter will maintain OSF/1 1.1 semantics.

### Execution Environment

The **lock\_init** kernel service can be called from the process environment only.

### Return Values

The **lock\_init** kernel service has no return values.

#### Related reference:

“lock\_alloc Kernel Service” on page 338

“lock\_free Kernel Service” on page 341

#### Related information:

Understanding Locking

Locking Kernel Services

## lock\_islocked Kernel Service

### Purpose

Tests whether a complex lock is locked.

### Syntax

```
#include <sys/lock_def.h>
```

```
int lock_islocked ( lock_addr)  
complex_lock_t lock_addr;
```

### Parameter

Item	Description
<i>lock_addr</i>	Specifies the address of the lock word to test.

## Description

The **lock\_islocked** kernel service determines whether the specified complex lock is free, or is locked in either shared-read or exclusive-write mode.

## Execution Environment

The **lock\_islocked** kernel service can be called from the process environment only.

## Return Values

Item	Description
TRUE	Indicates that the lock was locked.
FALSE	Indicates that the lock was free.

### Related reference:

“lock\_init Kernel Service” on page 342

### Related information:

Understanding Locking

Locking Kernel Services

## lockl Kernel Service

### Purpose

Locks a conventional process lock.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/lockl.h>
```

```
int lockl ( lock_word, flags)
lock_t *lock_word;
int flags;
```

### Parameters

Item	Description
<i>lock_word</i>	Specifies the address of the lock word.

**Item**  
*flags*

**Description**

Specifies the flags that control waiting for a lock. The *flags* parameter is used to control how signals affect waiting for a lock. The four flags are:

**LOCK\_NDELAY**

Controls whether the caller waits for the lock. Setting the flag causes the request to be terminated. The lock is assigned to the caller. Not setting the flag causes the caller to wait until the lock is not owned by another process before the lock is assigned to the caller.

**LOCK\_SHORT**

Prevents signals from terminating the wait for the lock. **LOCK\_SHORT** is the default flag for the **lockl** Kernel Service. This flag causes non-preemptive sleep.

**LOCK\_SIGRET**

Causes the wait for the lock to be terminated by an unmasked signal.

**LOCK\_SIGWAKE**

Causes the wait for the lock to be terminated by an unmasked signal and control transferred to the return from the last operation by the **setjmpx** kernel service.

**Note:** The **LOCK\_SIGRET** flag overrides the **LOCK\_SIGWAKE** flag.

## Description

**Note:** The **lockl** kernel service is provided for compatibility only and should not be used in new code, which should instead use simple locks or complex locks.

The **lockl** kernel service locks a conventional lock

The lock word can be located in shared memory. It must be in the process's address space when the **lockl** or **unlockl** services are called. The kernel accesses the lock word only while executing under the caller's process.

The *lock\_word* parameter is typically part of the data structure that describes the resource managed by the lock. This parameter must be initialized to the **LOCK\_AVAIL** value before the first call to the **lockl** service. Only the **lockl** and **unlockl** services can alter this parameter while the lock is in use.

The **lockl** service is nestable. The caller should use the **LOCK\_SUCC** value for determining when to call the **unlockl** service to unlock the conventional lock.

The **lockl** service temporarily assigns the owner the process priority of the most favored waiter for the lock.

A process must release all locks before terminating or leaving kernel mode. Signals are not delivered to kernel processes while those processes own any lock. "Understanding System Call Execution" in *Kernel Extensions and Device Support Programming Concepts* discusses how system calls can use the **lockl** service when accessing global data.

## Execution Environment

The **lockl** kernel service can be called from the process environment only.

## Return Values



Item	Description
LOCK_SUCC	Indicates that the process does not already own the lock or the lock is not owned by another process when the <i>flags</i> parameter is set to <b>LOCK_NDELAY</b> .
LOCK_NEST	Indicates that the process already owns the lock or the lock is not owned by another process when the <i>flags</i> parameter is set to <b>LOCK_NDELAY</b> .
LOCK_FAIL	Indicates that the lock is owned by another process when the <i>flags</i> parameter is set to <b>LOCK_NDELAY</b> .
LOCK_SIG	Indicates that the wait is terminated by a signal when the <i>flags</i> parameter is set to <b>LOCK_SIGRET</b> .

**Related reference:**

“unlockl Kernel Service” on page 519

**Related information:**

Understanding Locking

Locking Kernel Services

## lock\_mine Kernel Service

### Purpose

Checks whether a simple or complex lock is owned by the caller.

### Syntax

```
#include <sys/lock_def.h>
```

```
boolean_t lock_mine ( lock_addr)
void *lock_addr;
```

### Parameter

Item	Description
<i>lock_addr</i>	Specifies the address of the lock word to check.

### Description

The **lock\_mine** kernel service checks whether the specified simple or complex lock is owned by the calling kernel thread. Because a complex lock held in shared-read mode has no owner, the service returns **FALSE** in this case. This kernel service is provided to assist with debugging.

### Execution Environment

The **lock\_mine** kernel service can be called from the process environment only.

### Return Values

Item	Description
<b>TRUE</b>	Indicates that the calling kernel thread owns the lock.
<b>FALSE</b>	Indicates that the calling kernel thread does not own the lock, or that a complex lock is held in shared-read mode.

**Related reference:**

“lock\_init Kernel Service” on page 342

“lock\_read or lock\_try\_read Kernel Service” on page 346

“lock\_write or lock\_try\_write Kernel Service” on page 348

**Related information:**

Locking Kernel Services

## lock\_read or lock\_try\_read Kernel Service Purpose

Locks a complex lock in shared-read mode.

### Syntax

```
#include <sys/lock_def.h>
```

```
void lock_read ( lock_addr)  
complex_lock_t lock_addr;
```

```
boolean_t lock_try_read ( lock_addr)  
complex_lock_t lock_addr;
```

### Parameter

Item	Description
<i>lock_addr</i>	Specifies the address of the lock word to lock.

### Description

The **lock\_read** kernel service locks the specified complex lock in shared-read mode; it blocks if the lock is locked in exclusive-write mode. The lock must previously have been initialized with the **lock\_init** kernel service. The **lock\_read** kernel service has no return values.

The **lock\_try\_read** kernel service tries to lock the specified complex lock in shared-read mode; it returns immediately if the lock is locked in exclusive-write mode, otherwise it locks the lock in shared-read mode. The lock must previously have been initialized with the **lock\_init** kernel service.

### Execution Environment

The **lock\_read** and **lock\_try\_read** kernel services can be called from the process environment only.

### Return Values

The **lock\_try\_read** kernel service has the following return values:

Item	Description
TRUE	Indicates that the lock was successfully acquired in shared-read mode.
FALSE	Indicates that the lock was not acquired.

### Related reference:

“lock\_init Kernel Service” on page 342

“lock\_islocked Kernel Service” on page 342

“lock\_done Kernel Service” on page 340

### Related information:

Understanding Locking

Locking Kernel Services

## lock\_read\_to\_write or lock\_try\_read\_to\_write Kernel Service Purpose

Upgrades a complex lock from shared-read mode to exclusive-write mode.

## Syntax

```
#include <sys/lock_def.h>
```

```
boolean_t lock_read_to_write ( lock_addr)  
complex_lock_t lock_addr;
```

```
boolean_t lock_try_read_to_write ( lock_addr)  
complex_lock_t lock_addr;
```

## Parameter

Item	Description
<i>lock_addr</i>	Specifies the address of the lock word to be converted from read-shared to write-exclusive mode.

## Description

The `lock_read_to_write` and `lock_try_read_to_write` kernel services try to upgrade the specified complex lock from shared-read mode to exclusive-write mode. The lock is successfully upgraded if no other thread has already requested write-exclusive access for this lock. If the lock cannot be upgraded, it is no longer held on return from the `lock_read_to_write` kernel service; it is still held in shared-read mode on return from the `lock_try_read_to_write` kernel service.

The calling kernel thread must hold the lock in shared-read mode.

## Execution Environment

The `lock_read_to_write` and `lock_try_read_to_write` kernel services can be called from the process environment only.

## Return Values

The following only apply to `lock_read_to_write`:

Item	Description
TRUE	Indicates that the lock was not upgraded and is no longer held.
FALSE	Indicates that the lock was successfully upgraded to exclusive-write mode.

The following only apply to `lock_try_read_to_write`:

Item	Description
TRUE	Indicates that the lock was successfully upgraded to exclusive-write mode.
FALSE	Indicates that the lock was not upgraded and is held in read mode.

## Related reference:

“lock\_init Kernel Service” on page 342

“lock\_islocked Kernel Service” on page 342

“lock\_done Kernel Service” on page 340

## Related information:

Understanding Locking

Locking Kernel Services

## **lock\_set\_recursive Kernel Service**

### **Purpose**

Prepares a complex lock for recursive use.

### **Syntax**

```
#include <sys/lock_def.h>
```

```
void lock_set_recursive ( lock_addr)  
complex_lock_t lock_addr;
```

### **Parameter**

Item	Description
<i>lock_addr</i>	Specifies the address of the lock word to be prepared for recursive use.

### **Description**

The **lock\_set\_recursive** kernel service prepares the specified complex lock for recursive use. A complex lock cannot be nested until the **lock\_set\_recursive** kernel service is called for it. The calling kernel thread must hold the specified complex lock in write-exclusive mode.

When a complex lock is used recursively, the **lock\_done** kernel service must be called once for each time that the thread is locked in order to unlock the lock.

Only the kernel thread which calls the **lock\_set\_recursive** kernel service for a lock may acquire that lock recursively.

### **Execution Environment**

The **lock\_set\_recursive** kernel service can be called from process environment only.

### **Return Values**

The **lock\_set\_recursive** kernel service has no return values.

#### **Related reference:**

“lock\_init Kernel Service” on page 342

“lock\_done Kernel Service” on page 340

“lock\_write or lock\_try\_write Kernel Service”

“lock\_clear\_recursive Kernel Service” on page 339

#### **Related information:**

Understanding Locking

Locking Kernel Services

## **lock\_write or lock\_try\_write Kernel Service**

### **Purpose**

Locks a complex lock in exclusive-write mode.

### **Syntax**

```
#include <sys/lock_def.h>
```

```
void lock_write ( lock_addr)
complex_lock_t lock_addr;
```

```
boolean_t lock_try_write ( lock_addr)
complex_lock_t lock_addr;
```

## Parameter

Item	Description
<i>lock_addr</i>	Specifies the address of the lock word to lock.

## Description

The **lock\_write** kernel service locks the specified complex lock in exclusive-write mode; it blocks if the lock is busy. The lock must have been previously initialized with the **lock\_init** kernel service. The **lock\_write** kernel service has no return values.

The **lock\_try\_write** kernel service tries to lock the specified complex lock in exclusive-write mode; it returns immediately without blocking if the lock is busy. The lock must have been previously initialized with the **lock\_init** kernel service.

## Execution Environment

The **lock\_write** and **lock\_try\_write** kernel services can be called from the process environment only.

## Return Values

The **lock\_try\_write** kernel service has the following parameters:

Item	Description
TRUE	Indicates that the lock was successfully acquired.
FALSE	Indicates that the lock was not acquired.

### Related reference:

“lock\_done Kernel Service” on page 340

“lock\_read\_to\_write or lock\_try\_read\_to\_write Kernel Service” on page 346

### Related information:

Understanding Locking

Locking Kernel Services

## lock\_write\_to\_read Kernel Service

### Purpose

Downgrades a complex lock from exclusive-write mode to shared-read mode.

### Syntax

```
#include <sys/lock_def.h>
```

```
void lock_write_to_read ( lock_addr)
complex_lock_t lock_addr;
```

### Parameter

Item	Description
<i>lock_addr</i>	Specifies the address of the lock word to be downgraded from exclusive-write to shared-read mode.

## Description

The **lock\_write\_to\_read** kernel service downgrades the specified complex lock from exclusive-write mode to shared-read mode. The calling kernel thread must hold the lock in exclusive-write mode.

Once the lock has been downgraded to shared-read mode, other kernel threads will also be able to acquire it in shared-read mode.

## Execution Environment

The **lock\_write\_to\_read** kernel service can be called from the process environment only.

## Return Values

The **lock\_write\_to\_read** kernel service has no return values.

### Related reference:

“lock\_islocked Kernel Service” on page 342

“lock\_read\_to\_write or lock\_try\_read\_to\_write Kernel Service” on page 346

### Related information:

Understanding Locking

## loifp Kernel Service

### Purpose

Returns the address of the software loopback interface structure.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
struct ifnet *loifp ()
```

### Description

The **loifp** kernel service returns the address of the **ifnet** structure associated with the software loopback interface. The interface address can be used to examine the interface flags. This address can also be used to determine whether the **looutput** kernel service can be called to send a packet through the loopback interface.

### Execution Environment

The **loifp** kernel service can be called from either the process or interrupt environment.

### Return Values

The **loifp** service returns the address of the **ifnet** structure describing the software loopback interface.

### Related reference:

“looutput Kernel Service” on page 353

### Related information:

Network Kernel Services

## longjmpx Kernel Service

### Purpose

Allows exception handling by causing execution to resume at the most recently saved context.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
int longjmpx ( ret_val)
int ret_val;
```

### Parameters

Item	Description
<i>ret_val</i>	Specifies the return value to be supplied on the return from the <b>setjmpx</b> kernel service for the resumed context. This value normally indicates the type of exception that has occurred.

### Description

The **longjmpx** kernel service causes the normal execution flow to be modified so that execution resumes at the most recently saved context. The kernel mode lock is reacquired if it is necessary. The interrupt priority level is reset to that of the saved context.

The **longjmpx** service internally calls the **clrjmpx** service to remove the jump buffer specified by the *jump\_buffer* parameter from the list of contexts to be resumed. The **longjmpx** service always returns a nonzero value when returning to the restored context. Therefore, if the value of the *ret\_val* parameter is 0, the **longjmpx** service returns an **EINTR** value to the restored context.

If there is no saved context to resume, the system crashes.

### Execution Environment

The **longjmpx** kernel service can be called from either the process or interrupt environment.

### Return Values

A successful call to the **longjmpx** service does not return to the caller. Instead, it causes execution to resume at the return from a previous **setjmpx** call with the return value of the *ret\_val* parameter.

#### Related reference:

“**clrjmpx** Kernel Service” on page 43

“**setjmpx** Kernel Service” on page 468

#### Related information:

Understanding Exception Handling

Process and Exception Management Kernel Services

## lookupvp, lookupname, lookupname\_cur Kernel Services

### Purpose

Retrieves the v-node that corresponds to the named path.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/uio.h>
```

```
int lookupvp ( namep, flags, compvpp, crp)
char *namep;
int flags;
struct vnode **compvpp;
struct ucred *crp;
```

```
int lookupname ( namep, seg, flags, dirvpp, compvpp, crp)
char *namep;
int seg;
int flags;
struct vnode **dirvpp;
struct vnode **compvpp;
struct cred *crp;
```

```
int lookupname_cur ( namep, seg, flags, dirvpp, compvpp, curdvp, crp)
char *namep;
int seg;
int flags;
struct vnode **dirvpp;
struct vnode **compvpp;
struct vnode **curdvp;
struct cred *crp;
```

## Parameters

Item	Description
<i>crp</i>	Points to the <b>cred</b> structure. This structure contains data that the file system can use to validate access permission.
<i>namep</i>	Points to a character string path name.
<i>flags</i>	Specifies lookup directives, including these six flags: <b>L_LOC</b> The path-name resolution must not cross a mount point into another file system implementation. <b>L_NOFOLLOW</b> If the final component of the path name resolves to a symbolic link, the link is not to be traversed. <b>L_NOXMOUNT</b> If the final component of the path name resolves to a mounted-over object, the mounted-over object, rather than the root of the next virtual file system, is to be returned. <b>L_CRT</b> The object is to be created. <b>L_DEL</b> The object is to be deleted. <b>L_EROFS</b> An error is to be returned if the object resides in a read-only file system.
<i>seg</i>	Specifies whether the <i>namep</i> buffer is in user space (UIO_USERSPACE) or kernel space (UIO_SYSSPACE).
<i>compvpp</i>	Points to the location where the vnode pointer for the named object is to be returned to the calling routine.
<i>dirvpp</i>	Points to the location where the vnode pointer for the directory containing the named object is to be returned.
<i>curdvp</i>	Points to the vnode for a current directory to be used instead of <i>u_cdir</i> .

## Description

The **lookupvp** kernel service provides translation of the path name provided by the *namep* parameter into a virtual file system node. The **lookupvp** service provides a flexible interface to path-name resolution by



regarding the *flags* parameter values as directives to the lookup process. The lookup process is a cooperative effort between the logical file system and underlying virtual file systems (VFS). Several v-node and VFS operations are employed to:

- Look up individual name components
- Read symbolic links
- Cross mount points

The **lookupvp** kernel service determines the process's current and root directories by consulting the `u_cdir` and `u_rdir` fields in the `u` structure. Information about the virtual file system and file system installation for transient v-nodes is obtained from each name component's `vfs` or `gfs` structure. The **lookupvp** kernel service assumes that the named path is in kernel address space.

The **lookupname** kernel service provides the same service as the **lookupvp** kernel service, but allows the caller to specify whether the path name is in kernel or user space. It also provides the ability to retrieve the vnode for the directory containing the named object. The **lookupname\_cur** kernel service further extends the interface by allowing the lookup to proceed relative to the given *curdvp* directory.

The vnodes returned by the **lookup** services are held. The calling routine is responsible for releasing the hold by calling the **vnop\_rele** entry point when it completes its operation.

## Execution Environment

The **lookup** kernel services can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
errno	Indicates an error. This number is defined in the <code>/usr/include/sys/errno.h</code> file.

## Related information:

Understanding Data Structures and Header Files for Virtual File Systems

Virtual File System Overview

Virtual File System (VFS) Kernel Services

## looutput Kernel Service

### Purpose

Sends data through a software loopback interface.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
int looutput ( ifp, m0, dst)
struct ifnet *ifp;
struct mbuf *m0;
struct sockaddr *dst;
```

### Parameters

Item	Description
<i>ifp</i>	Specifies the address of an <b>ifnet</b> structure describing the software loopback interface.
<i>m0</i>	Specifies an <b>mbuf</b> chain containing output data.
<i>dst</i>	Specifies the address of a <b>sockaddr</b> structure that specifies the destination for the data.

## Description

The **looutput** kernel service sends data through a software loopback interface. The data in the *m0* parameter is passed to the input handler of the protocol specified by the *dst* parameter.

## Execution Environment

The **looutput** kernel service can be called from either the process or interrupt environment.

## Return Values

Item	Description
0	Indicates that the data was successfully sent.
ENOBUFS	Indicates that resource allocation failed.
EAFNOSUPPORT	Indicates that the address family specified by the <i>dst</i> parameter is not supported.

### Related reference:

“loifp Kernel Service” on page 350

### Related information:

Network Kernel Services

## ltpin Kernel Service

### Purpose

Pins the address range in the system (kernel) space and frees the page space for the associated pages.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/pin.h>
int ltpin (addr, length)
caddr_t  addr;
int      length;
```

### Parameters

Item	Description
<i>addr</i>	Specifies the address of the first byte to pin.
<i>length</i>	Specifies the number of bytes to pin.

## Description

The **ltpin** (long term pin) kernel service pins the real memory pages touched by the address range specified by the *addr* and *length* parameters in the system (kernel) address space. It pins the real-memory pages to ensure that page faults do not occur for memory references in this address range. The **ltpin** kernel service increments the long-term pin count for each real-memory page. While either the long-term or short-term pin count is nonzero, the page cannot be paged out of real memory.

The **Itpin** kernel service pins either the entire address range or none of it. Only a limited number of pages are pinned in the system. If there are not enough unpinned pages in the system, the **Itpin** kernel service returns an error code. The **Itpin** kernel service is not a published interface.

**Note:** The operating system pins only whole pages at a time. Therefore, if the requested range is not aligned on a page boundary, then memory outside this range is also pinned.

The **Itpin** kernel service can only be called for addresses within the system (kernel) address space.

## Return Values

Item	Description
0	Indicates successful completion.
EINVAL	Indicates that the <i>length</i> parameter has a negative value. Otherwise, the area of memory beginning at the address of the first byte to pin (the <b>addr</b> parameter) and extending for the number of bytes specified by the <i>length</i> parameter is not defined.
EIO	Indicates that a permanent I/O error occurred while referencing data.
ENOMEM	Indicates that the <b>pin</b> kernel service was unable to pin due to insufficient real memory or exceeding the system-wide pin count.
ENOSPC	Indicates insufficient file system or paging space.

## Related reference:

“Itunpin Kernel Service”

## Itunpin Kernel Service

### Purpose

Unpins the address range in system (kernel) address space and reallocates paging space for the specified region.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/pin.h>

int  Itunpin (addr, length)
caddr_t  addr;
int      length;
```

### Parameters

Item	Description
<i>addr</i>	Specifies the address of the first byte to unpin.
<i>length</i>	Specifies the number of bytes to unpin.

### Description

The **Itunpin** kernel service decreases the long-term pin count of each page in the address range. When the long-term pin count becomes 0, the backing storage (paging space) for the memory region is allocated and assigned to the pages. When both the long-term and short-term pin counts are 0, the page is no longer pinned and the **Itunpin** kernel service will assert. If allocating backing pages would put the system below the low paging space threshold, the call waits until paging space becomes available.

The **Itunpin** kernel service can only be called with addresses in the system (kernel) address space from the process environment.

## Return Values

Item	Description
0	Indicates successful completion.
EINVAL	Indicates that the <i>length</i> parameter is a negative value.
EIO	Indicates that a permanent I/O error occurred while referencing data.

### Related reference:

“ltpin Kernel Service” on page 354

## m

The following kernel services begin with the with the letter m.

### m\_adj Kernel Service

#### Purpose

Adjusts the size of an **mbuf** chain.

#### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/mbuf.h>
```

```
void m_adj ( m, diff)
struct mbuf *m;
int diff;
```

#### Parameters

Item	Description
<i>m</i>	Specifies the <b>mbuf</b> chain to be adjusted.
<i>diff</i>	Specifies the number of bytes to be removed.

#### Description

The **m\_adj** kernel service adjusts the size of an **mbuf** chain by the number of bytes specified by the *diff* parameter. If the number specified by the *diff* parameter is nonnegative, the bytes are removed from the front of the chain. If this number is negative, the alteration is done from back to front.

#### Execution Environment

The **m\_adj** kernel service can be called from either the process or interrupt environment.

#### Return Values

The **m\_adj** service has no return values.

#### Related information:

I/O Kernel Services

### mbreq Structure for mbuf Kernel Services

#### Purpose

Contains **mbuf** structure registration information for the **m\_reg** and **m\_dereg** kernel services.

## Syntax

```
#include <sys/mbuf.h>
struct mbreq {
    int low_mbuf;
    int low_clust;
    int initial_mbuf;
    int initial_clust;
}
```

## Parameters

Item	Description
<i>low_mbuf</i>	Specifies the <b>mbuf</b> structure low-water mark.
<i>low_clust</i>	Specifies the page-sized <b>mbuf</b> structure low-water mark.
<i>initial_mbuf</i>	Specifies the initial allocation of <b>mbuf</b> structures.
<i>initial_clust</i>	Specifies the initial allocation of page-sized <b>mbuf</b> structures.

## Description

The **mbreq** structure specifies the **mbuf** structure usage expectations for a user of **mbuf** kernel services.

### Related reference:

“m\_dereg Kernel Service” on page 364

“m\_reg Kernel Service” on page 373

### Related information:

I/O Kernel Services

## **mbstat** Structure for mbuf Kernel Services Purpose

Contains **mbuf** usage statistics.

## Syntax

```
#include <sys/mbuf.h>
struct mbstat {
    ulong m_mbufs;
    ulong m_clusters;
    ulong m_spare;
    ulong m_clfree;
    ulong m_drops;
    ulong m_wait;
    ulong m_drain;
    short m_mtypes[256];
}
```

## Parameters

Item	Description
<i>m_mbufs</i>	Specifies the number of <b>mbuf</b> structures allocated.
<i>m_clusters</i>	Specifies the number of clusters allocated.
<i>m_spare</i>	Specifies the spare field.
<i>m_clfree</i>	Specifies the number of free clusters.
<i>m_drops</i>	Specifies the times failed to find space.
<i>m_wait</i>	Specifies the times waited for space.
<i>m_drain</i>	Specifies the times drained protocols for space.
<i>m_mtypes</i>	Specifies the type-specific <b>mbuf</b> structure allocations.

## Description

The **mbstat** structure provides usage information for the **mbuf** services. Statistics can be viewed through the **netstat -m** command.

### Related information:

netstat subroutine

I/O Kernel Services

## m\_cat Kernel Service

### Purpose

Appends one **mbuf** chain to the end of another.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/mbuf.h>
```

```
void m_cat ( m, n)
struct mbuf *m;
struct mbuf *n;
```

### Parameters

Item	Description
<i>m</i>	Specifies the <b>mbuf</b> chain to be appended to.
<i>n</i>	Specifies the <b>mbuf</b> chain to append.

## Description

The **m\_cat** kernel service appends an **mbuf** chain specified by the *n* parameter to the end of **mbuf** chain specified by the *m* parameter. Where possible, compaction is performed.

## Execution Environment

The **m\_cat** kernel service can be called from either the process or interrupt environment.

## Return Values

The **m\_cat** service has no return values.

### Related information:

I/O Kernel Services

## m\_clattach Kernel Service

### Purpose

Allocates an **mbuf** structure and attaches an external cluster.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/mbuf.h>
```

```
struct mbuf *
m_clattach( ext_buf, ext_free, ext_size, ext_arg, wait)
caddr_t ext_buf;
int (*ext_free)();
int ext_size;
int ext_arg;
int wait;
```

### Parameters

Item	Description
<i>ext_buf</i>	Specifies the address of the external data area.
<i>ext_free</i>	Specifies the address of a function to be called when this <b>mbuf</b> structure is freed.
<i>ext_size</i>	Specifies the length of the external data area.
<i>ext_arg</i>	Specifies an argument to pass to the above function.
<i>wait</i>	Specifies either the <b>M_WAIT</b> or <b>M_DONTWAIT</b> value.

### Description

The **m\_clattach** kernel service allocates an **mbuf** structure and attaches the cluster specified by the *ext\_buf* parameter. This data is owned by the caller. The *m\_data* field of the returned **mbuf** structure points to the caller's data. Interrupt handlers can call this service only with the *wait* parameter set to **M\_DONTWAIT**.

**Note:** The **m\_clattach** kernel service replaces the **m\_clgetx** kernel service, which is no longer supported.

The calling function is required to fill out the **mbuf** structure sufficiently to support normal usage. This includes support for the DMA functions during network transmission. To support DMA functions, the **ext\_hasxm** flag field needs to be set to true and the **ext\_xmemd** structure needs to be filled out. For buffers allocated from the kernel pinned heap, the **ext\_xmemd.aspace\_id** field should be set to **XMEM\_GLOBAL**.

### Execution Environment

The **m\_clattach** kernel service can be called from either the process or interrupt environment.

### Return Values

The **m\_clattach** kernel service returns the address of an allocated **mbuf** structure. If the *wait* parameter is set to **M\_DONTWAIT** and there are no free **mbuf** structures, the **m\_clattach** service returns null.

### Related information:

I/O Kernel Services

## m\_clget Macro for mbuf Kernel Services

### Purpose

Allocates a page-sized **mbuf** structure cluster.

### Syntax

```
#include <sys/mbuf.h>
```

```
int m_clget ( m )  
struct mbuf *m;
```

### Parameter

Item	Description
<i>m</i>	Specifies the <b>mbuf</b> structure with which the cluster is to be associated.

### Description

The **m\_clget** macro allocates a page-sized **mbuf** cluster and attaches it to the given **mbuf** structure. If successful, the length of the **mbuf** structure is set to **CLBYTES**.

### Execution Environment

The **m\_clget** macro can be called from either the process or interrupt environment.

### Return Values

Item	Description
1	Indicates successful completion.
0	Indicates an error.

### Related reference:

“m\_clgetm Kernel Service”

### Related information:

I/O Kernel Services

## m\_clgetm Kernel Service

### Purpose

Allocates and attaches an external buffer.

### Syntax

```
#include <sys/types.h>  
#include <sys/errno.h>  
#include <sys/mbuf.h>  
#include <net/net_globals.h>
```

```
int  
m_clgetm( m, how, size )  
struct mbuf *m;  
int how;  
int size;
```



## Parameters

Item	Description
<i>m</i>	Specifies the <b>mbuf</b> structure that the cluster will be associated with.
<i>how</i>	Specifies either the <b>M_DONTWAIT</b> or <b>M_WAIT</b> value.
<i>size</i>	Specifies the size of external cluster to attach. Any value less than <b>MAXALLOCSAVE</b> is valid. For larger values, <b>M_WAIT</b> must be specified.

## Description

The **m\_clgetm** service allocates an **mbuf** cluster of the specified number of bytes and attaches it to the **mbuf** structure indicated by the *m* parameter. If successful, the **m\_clgetm** service sets the **M\_EXT** flag.

## Execution Environment

The **m\_clgetm** kernel service can be called from either the process or interrupt environment.

An interrupt handler can specify the *wait* parameter as **M\_DONTWAIT** only.

## Return Values

Item	Description
1	Indicates a successful operation.

If there are no free **mbuf** structures, the **m\_clgetm** kernel service returns a null value.

### Related reference:

“**m\_freem** Kernel Service” on page 366

“**m\_get** Kernel Service” on page 367

“**m\_clget** Macro for **mbuf** Kernel Services” on page 360

### Related information:

I/O Kernel Services

## **m\_collapse** Kernel Service

### Purpose

Guarantees that an **mbuf** chain contains no more than a given number of **mbuf** structures.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/mbuf.h>
```

```
struct mbuf *m_collapse ( m, size)
struct mbuf *m;
int size;
```

## Parameters

Item	Description
<i>m</i>	Specifies the <b>mbuf</b> chain to be collapsed.
<i>size</i>	Denotes the maximum number of <b>mbuf</b> structures allowed in the chain.

## Description

The **m\_collapse** kernel service reduces the number of **mbuf** structures in an **mbuf** chain to the number of **mbuf** structures specified by the *size* parameter. The **m\_collapse** service accomplishes this by copying data into page-sized **mbuf** structures until the chain is of the desired length. (If required, more than one page-sized **mbuf** structure is used.)

## Execution Environment

The **m\_collapse** kernel service can be called from either the process or interrupt environment.

## Return Values

If the chain cannot be collapsed into the number of **mbuf** structures specified by the *size* parameter, a value of null is returned and the original chain is deallocated. Upon successful completion, the head of the altered **mbuf** chain is returned.

### Related information:

I/O Kernel Services

## **m\_copy** Macro for mbuf Kernel Services Purpose

Creates a copy of all or part of a list of **mbuf** structures.

## Syntax

```
#include <sys/mbuf.h>
```

```
struct mbuf *m_copy ( m, off, len)
struct mbuf *m;
int off;
int len;
```

## Parameters

Item	Description
<i>m</i>	Specifies the <b>mbuf</b> structure, or the head of a list of <b>mbuf</b> structures, to be copied.
<i>off</i>	Specifies an offset into data from which copying starts.
<i>len</i>	Denotes the total number of bytes to copy.

## Description

The **m\_copy** macro makes a copy of the structure specified by the *m* parameter. The copy begins at the specified bytes (represented by the *off* parameter) and continues for the number of bytes specified by the *len* parameter. If the *len* parameter is set to **M\_COPYALL**, the entire **mbuf** chain is copied.

## Execution Environment

The **m\_copy** macro can be called from either the process or interrupt environment.

## Return Values

Upon successful completion, the address of the copied list (the **mbuf** structure that heads the list) is returned. If the copy fails, a value of null is returned.

### Related reference:

“m\_copydata Kernel Service”

“m\_copym Kernel Service” on page 364

### Related information:

I/O Kernel Services

## m\_copydata Kernel Service

### Purpose

Copies data from an **mbuf** chain to a specified buffer.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/mbuf.h>
```

```
void m_copydata (m, off, len, cp)
struct mbuf * m;
int off;
int len;
caddr_t cp;
```

### Parameters

Item	Description
<i>m</i>	Indicates the <b>mbuf</b> structure, or the head of a list of <b>mbuf</b> structures, to be copied.
<i>off</i>	Specifies an offset into data from which copying starts.
<i>len</i>	Denotes the total number of bytes to copy.
<i>cp</i>	Points to a data buffer into which to copy the <b>mbuf</b> data.

### Description

The **m\_copydata** kernel service makes a copy of the structure specified by the *m* parameter. The copy begins at the specified bytes (represented by the *off* parameter) and continues for the number of bytes specified by the *len* parameter. The data is copied into the buffer specified by the *cp* parameter.

### Execution Environment

The **m\_copydata** kernel service can be called from either the process or interrupt environment.

### Return Values

The **mcopydata** service has no return values.

### Related reference:

“m\_copy Macro for mbuf Kernel Services” on page 362

### Related information:

I/O Kernel Services

## **m\_copym Kernel Service**

### **Purpose**

Creates a copy of all or part of a list of **mbuf** structures.

### **Syntax**

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/mbuf.h>
```

```
struct mbuf *
m_copym( m, off, len, wait)
struct mbuf m;
int off;
int len;
int wait;
```

### **Parameters**

Item	Description
<i>m</i>	Specifies the <b>mbuf</b> structure to be copied.
<i>off</i>	Specifies an offset into data from which copying will start.
<i>len</i>	Specifies the total number of bytes to copy.
<i>wait</i>	Specifies either the <b>M_DONTWAIT</b> or <b>M_WAIT</b> value.

### **Description**

The **m\_copym** kernel service makes a copy of the **mbuf** structure specified by the *m* parameter starting at the specified offset from the beginning and continuing for the number of bytes specified by the *len* parameter. If the *len* parameter is set to **M\_COPYALL**, the entire **mbuf** chain is copied.

If the **mbuf** structure specified by the *m* parameter has an external buffer attached (that is, the **M\_EXT** flag is set), the copy is done by reference to the external cluster. In this case, the data must not be altered or both copies will be changed. Interrupt handlers can specify the *wait* parameter as **M\_DONTWAIT** only.

### **Execution Environment**

The **m\_copym** kernel service can be called from either the process or interrupt environment.

### **Return Values**

The address of the copy is returned upon successful completion. If the copy fails, null is returned. If the *wait* parameter is set to **M\_DONTWAIT** and there are no free **mbuf** structures, the **m\_copym** kernel service returns a null value.

#### **Related reference:**

“m\_copydata Kernel Service” on page 363

“m\_copy Macro for mbuf Kernel Services” on page 362

#### **Related information:**

I/O Kernel Services

## **m\_dereg Kernel Service**

### **Purpose**

Deregisters expected **mbuf** structure usage.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/mbuf.h>
```

```
void m_dereg ( mbp)
struct mbreq mbp;
```

## Parameter

Item	Description
<i>mbp</i>	Defines the address of an <b>mbreq</b> structure that specifies expected <b>mbuf</b> usage.

## Description

The **m\_dereg** kernel service deregisters requirements previously registered with the **m\_reg** kernel service. The **m\_dereg** service is mandatory if the **m\_reg** service is called.

## Execution Environment

The **m\_dereg** kernel service can be called from the process environment only.

## Return Values

The **m\_dereg** service has no return values.

### Related reference:

“mbreq Structure for mbuf Kernel Services” on page 356

“m\_reg Kernel Service” on page 373

### Related information:

I/O Kernel Services

## m\_free Kernel Service

### Purpose

Frees an **mbuf** structure and any associated external storage area.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/mbuf.h>
```

```
struct mbuf *m_free( m)
struct mbuf *m;
```

## Parameter

Item	Description
<i>m</i>	Specifies the <b>mbuf</b> structure to be freed.

## Description

The **m\_free** kernel service returns an **mbuf** structure to the buffer pool. If the **mbuf** structure specified by the *m* parameter has an attached cluster (that is, a paged-size **mbuf** structure), the **m\_free** kernel service also frees the associated external storage.

## Execution Environment

The **m\_free** kernel service can be called from either the process or interrupt environment.

## Return Values

If the **mbuf** structure specified by the *m* parameter is the head of an **mbuf** chain, the **m\_free** service returns the next **mbuf** structure in the chain. A null value is returned if the structure specified by the *m* parameter is not part of an **mbuf** chain.

### Related reference:

“m\_get Kernel Service” on page 367

### Related information:

I/O Kernel Services

## m\_freem Kernel Service

### Purpose

Frees an entire **mbuf** chain.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/mbuf.h>
```

```
void m_freem ( m )
struct mbuf *m;
```

### Parameter

Item	Description
<i>m</i>	Indicates the head of the <b>mbuf</b> chain to be freed.

## Description

The **m\_freem** kernel service starts the **m\_free** kernel service for each **mbuf** structure in the chain headed by the head specified by the *m* parameter.

## Execution Environment

The **m\_freem** kernel service can be called from either the process or interrupt environment.

## Return Values

The **m\_freem** service has no return values.

### Related reference:

“m\_free Kernel Service” on page 365

“m\_get Kernel Service”

**Related information:**

I/O Kernel Services

## m\_get Kernel Service

### Purpose

Allocates a memory buffer (mbuf) from the **mbuf** pool.

### Syntax

```
#include <sys/types.h>
```

```
#include <sys/errno.h>
```

```
#include <sys/mbuf.h>
```

```
struct mbuf *m_get ( wait, type)
```

```
int wait;
```

```
int type;
```

### Parameters

Item	Description
------	-------------

<i>wait</i>	Indicates the action to be taken if there are no free <b>mbuf</b> structures. Possible values are:
-------------	--

**M\_DONTWAIT**

Called from either an interrupt or process environment.

**M\_WAIT**

Called from a process environment.

<i>type</i>	Specifies a valid <b>mbuf</b> type, as listed in the <code>/usr/include/sys/mbuf.h</code> file.
-------------	---

### Description

The **m\_get** kernel service allocates an **mbuf** structure of the specified type. If the buffer pool is empty and the *wait* parameter is set to **M\_WAIT**, the **m\_get** kernel service does not return until an **mbuf** structure is available.

### Execution Environment

The **m\_get** kernel service can be called from either the process or interrupt environment.

An interrupt handler can specify the *wait* parameter as **M\_DONTWAIT** only.

### Return Values

Upon successful completion, the **m\_get** service returns the address of an allocated **mbuf** structure. If the *wait* parameter is set to **M\_DONTWAIT** and there are no free **mbuf** structures, the **m\_get** kernel service returns a null value.

**Related reference:**

“m\_free Kernel Service” on page 365

“m\_freem Kernel Service” on page 366

**Related information:**

I/O Kernel Services

## **m\_getclr Kernel Service**

### **Purpose**

Allocates and zeroes a memory buffer from the **mbuf** pool.

### **Syntax**

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/mbuf.h>
```

```
struct mbuf *m_getclr ( wait, type)
int wait;
int type;
```

### **Parameters**

Item	Description
<i>wait</i>	This flag indicates the action to be taken if there are no free <b>mbuf</b> structures. Possible values are:  <b>M_DONTWAIT</b> Called from either an interrupt or process environment.  <b>M_WAIT</b> Called from a process environment only.
<i>type</i>	Specifies a valid <b>mbuf</b> type, as listed in the <code>/usr/include/sys/mbuf.h</code> file.

### **Description**

The **m\_getclr** kernel service allocates an **mbuf** structure of the specified type. If the buffer pool is empty and the *wait* parameter is set to **M\_WAIT** value, the **m\_getclr** service does not return until an **mbuf** structure is available.

The **m\_getclr** kernel service differs from the **m\_get** kernel service in that the **m\_getclr** service zeroes the data portion of the allocated **mbuf** structure.

### **Execution Environment**

The **m\_getclr** kernel service can be called from either the process or interrupt environment. Interrupt handlers can call the **m\_getclr** service only with the *wait* parameter set to the **M\_DONTWAIT** value.

### **Return Values**

The **m\_getclr** kernel service returns the address of an allocated **mbuf** structure. If the *wait* parameter is set to the **M\_DONTWAIT** value and there are no free **mbuf** structures, the **m\_getclr** kernel service returns a null value.

#### **Related reference:**

- “m\_free Kernel Service” on page 365
- “m\_freem Kernel Service” on page 366
- “m\_get Kernel Service” on page 367

#### **Related information:**

I/O Kernel Services

## **m\_getclust Macro for mbuf Kernel Services**

### **Purpose**

Allocates an **mbuf** structure from the **mbuf** buffer pool and attaches a page-sized cluster.



## Syntax

```
#include <sys/mbuf.h>
```

```
struct mbuf *m_getclust ( wait, type)
int wait;
int type;
```

## Parameters

Item	Description
<i>wait</i>	Indicates the action to be taken if there are no available <b>mbuf</b> structures. Possible values are:  M_DONTWAIT Called from either an interrupt or process environment.  M_WAIT Called from a process environment only.
<i>type</i>	Specifies a valid <b>mbuf</b> type from the <code>/usr/include/sys/mbuf.h</code> file.

## Description

The **m\_getclust** macro allocates an **mbuf** structure of the specified type. If the allocation succeeds, the **m\_getclust** macro then attempts to attach a page-sized cluster to the structure.

If the buffer pool is empty and the *wait* parameter is set to **M\_WAIT**, the **m\_getclust** macro does not return until an **mbuf** structure is available.

## Execution Environment

The **m\_getclust** macro can be called from either the process or interrupt environment.

## Return Values

The address of an allocated **mbuf** structure is returned on success. If the *wait* parameter is set to **M\_DONTWAIT** and there are no free **mbuf** structures, the **m\_getclust** macro returns a null value.

### Related reference:

“m\_getclustm Kernel Service”

### Related information:

I/O Kernel Services

## m\_getclustm Kernel Service

### Purpose

Allocates an **mbuf** structure and attaches a cluster of the specified size, both from the **mbuf** buffer pool.

## Syntax

```
#include <sys/mbuf.h>
#include <net/net_globals.h>
```

```
struct mbuf *
m_getclustm( wait, type, size)
int wait;
int type;
int size;
```

## Parameters

Item	Description
<i>wait</i>	Specifies either the <code>M_DONTWAIT</code> or <code>M_WAIT</code> value.
<i>type</i>	Specifies a valid <code>mbuf</code> type from the <code>/usr/include/sys/mbuf.h</code> file.
<i>size</i>	Specifies the size of the external cluster to attach. Any value less than <code>MAXALLOCSAVE</code> is valid. For larger values, <code>M_WAIT</code> must be specified.

## Description

The `m_getclustm` service allocates an `mbuf` structure of the specified type. If successful, the `m_getclustm` service then attempts to attach a cluster of the indicated size (specified by the *size* parameter) to the `mbuf` structure. If the buffer pool is empty and the *wait* parameter is set to `M_WAIT`, the `m_get` service does not return until an `mbuf` structure is available. Interrupt handlers should call this service only with the *wait* parameter set to `M_DONTWAIT`.

## Execution Environment

The `m_getclustm` kernel service can be called from either the process or interrupt environment.

An interrupt handler can specify the *wait* parameter as `M_DONTWAIT` only.

## Return Values

The `m_getclustm` kernel service returns the address of an allocated `mbuf` structure on success. If the *wait* parameter is set to `M_DONTWAIT` and there are no free `mbuf` structures, the `m_getclustm` kernel service returns null.

### Related reference:

“`m_clget` Macro for `mbuf` Kernel Services” on page 360

“`m_free` Kernel Service” on page 365

“`m_freem` Kernel Service” on page 366

“`m_get` Kernel Service” on page 367

“`m_getclust` Macro for `mbuf` Kernel Services” on page 368

### Related information:

I/O Kernel Services

## `m_gethdr` Kernel Service

### Purpose

Allocates a header memory buffer from the `mbuf` pool.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/mbuf.h>
```

```
struct mbuf *
m_gethdr ( wait, type)
int wait;
int type;
```

### Parameters

Item	Description
<i>wait</i>	Specifies either the <code>M_DONTWAIT</code> or <code>M_WAIT</code> value.
<i>type</i>	Specifies the valid <code>mbuf</code> type from the <code>/usr/include/sys/mbuf.h</code> file.

## Description

The `m_gethdr` kernel service allocates an `mbuf` structure of the specified type. If the buffer pool is empty and the *wait* parameter is set to `M_WAIT`, the `m_gethdr` kernel service will not return until an `mbuf` structure is available. Interrupt handlers should call this kernel service only with the *wait* parameter set to `M_DONTWAIT`. The `M_PKTHDR` flag is set for the returned `mbuf` structure.

## Execution Environment

The `m_gethdr` kernel service can be called from either the process or interrupt environment.

An interrupt handler can specify the *wait* parameter as `M_DONTWAIT` only.

## Return Values

The address of an allocated `mbuf` structure is returned on success. If the *wait* parameter is set to `M_DONTWAIT` and there are no free `mbuf` structure, the `m_gethdr` kernel service returns null.

## Related Information

The `m_free` kernel service, `m_freem` kernel service.

I/O Kernel Services in *Kernel Extensions and Device Support Programming Concepts*.

### Related reference:

“`m_free` Kernel Service” on page 365

“`m_freem` Kernel Service” on page 366

### Related information:

I/O Kernel Services

## M\_HASCL Macro for mbuf Kernel Services Purpose

Determines if an `mbuf` structure has an attached cluster.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/mbuf.h>
```

```
struct mbuf * m;
M_HASCL (m);
```

## Parameter

Item	Description
<i>m</i>	Indicates the address of the <b>mbuf</b> structure in question.

## Description

The **M\_HASCL** macro determines if an **mbuf** structure has an attached cluster.

## Execution Environment

The **M\_HASCL** macro can be called from either the process or interrupt environment.

## Example

The **M\_HASCL** macro can be used as in the following example:

```
struct mbuf *m;
if (M_HASCL(m))
    printf("mbuf has attached cluster");
```

## Related information:

I/O Kernel Services

## **m\_pullup** Kernel Service

### Purpose

Adjusts an **mbuf** chain so that a given number of bytes is in contiguous memory in the data area of the head **mbuf** structure.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/mbuf.h>
```

```
struct mbuf *m_pullup ( m, size)
struct mbuf *m;
int size;
```

### Parameters

Item	Description
<i>m</i>	Specifies the <b>mbuf</b> chain to be adjusted.
<i>size</i>	Specifies the number of bytes to be contiguous.

### Description

The **m\_pullup** kernel service guarantees that the **mbuf** structure at the head of a chain has in contiguous memory within its data area at least the number of data bytes specified by the *size* parameter.

### Execution Environment

The **m\_pullup** kernel service can be called from either the process or interrupt environment.

### Return Values

Upon successful completion, the head structure in the altered **mbuf** chain is returned.

A value of null is returned and the original chain is deallocated under the following circumstances:

- The size of the chain is less than indicated by the *size* parameter.
- The number indicated by the *size* parameter is greater than the data portion of the head-size **mbuf** structure.

**Related information:**

I/O Kernel Services

## **m\_reg Kernel Service**

### **Purpose**

Registers expected **mbuf** usage.

### **Syntax**

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/mbuf.h>
```

```
void m_reg ( mbp)
struct mbreq mbp;
```

### **Parameter**

Item	Description
<i>mbp</i>	Defines the address of an <b>mbreq</b> structure that specifies expected <b>mbuf</b> usage.

### **Description**

The **m\_reg** kernel service lets users of **mbuf** services specify initial requirements. The **m\_reg** kernel service also allows the buffer pool low-water and deallocation marks to be adjusted based on expected usage. Its use is recommended for better control of the buffer pool.

When the number of free **mbuf** structures falls below the low-water mark, the total **mbuf** pool is expanded. When the number of free **mbuf** structures rises above the deallocation mark, the total **mbuf** pool is contracted and resources are returned to the system.

### **Execution Environment**

The **m\_reg** kernel service can be called from the process environment only.

### **Return Values**

The **m\_reg** service has no return values.

**Related reference:**

“mbreq Structure for mbuf Kernel Services” on page 356

“m\_dereg Kernel Service” on page 364

**Related information:**

I/O Kernel Services

## **md\_restart\_block\_read Kernel Service**

### **Purpose**

A copy of the **RESTART\_BLOCK** structure in the **NVRAM** header will be placed in the caller's buffer.

### **Syntax**

```
#include <sys/mdio.h>
```

```
int md_restart_block_read (md)
    struct mdio *md;
```

## Parameters

### Item Description

*md* Specifies the address of the **mdio** structure. The **mdio** structure contains the following fields:

<b>md_data</b>	Pointer to the data buffer.
<b>md_size</b>	Number of bytes in the data buffer.
<b>md_addr</b>	Contains the value PMMode on return in the least significant byte.

## Description

The RestartBlock which is in the **NVRAM** header will be copied to the user supplied buffer. This block is a communication vehicle for the software and the firmware.

## Return Values

Returns 0 for successful completion.

### Item Description

ENOMEM	Indicates that there was not enough room in the user supplied buffer to contain the RestartBlock.
EINVAL	Indicates this is not a PowerPC® reference platform.

## Prerequisite Information

Kernel Extensions and Device Driver Management Kernel Services in Kernel Extensions and Device Support Programming Concepts.

### Related information:

Machine Device Driver

## md\_restart\_block\_upd Kernel Service

### Purpose

The caller supplied RestartBlock will be copied to the **NVRAM** header.

## Syntax

```
#include <sys/mdio.h>
int md_restart_block_upd (md, pmmode)
    struct mdio *md;
    unsigned char pmmode;
```

## Description

The 8-bit value in *pmmode* will be stored into the **NVRAM** header at the PMMode offset. The RestartBlock which is in the caller's buffer will be copied to the **NVRAM** after the RestartBlock checksum is calculated and a new Crc1 value is computed.

## Parameters

Item	Description
<i>md</i>	Specifies the address of the <b>mdio</b> structure. The <b>mdio</b> structure contains the following fields:
	<b>md_data</b> Pointer to the RestartBlock structure..
<i>pmmode</i>	Value to be stored into PMMode in the NVRAM header.

## Return Values

Returns 0 for successful completion.

Item	Description
EINVAL	Indicates this is not a PowerPC reference platform.

## Prerequisite Information

Kernel Extensions and Device Driver Management Kernel Services in Kernel Extensions and Device Support Programming Concepts.

### Related information:

Machine Device Driver

## MTOCL Macro for mbuf Kernel Services Purpose

Converts a pointer to an **mbuf** structure to a pointer to the head of an attached cluster.

## Syntax

```
#include <sys/mbuf.h>
```

```
struct mbuf * m;  
MTOCL (m);
```

## Parameter

Item	Description
<i>m</i>	Indicates the address of the <b>mbuf</b> structure in question.

## Description

The **MTOCL** macro converts a pointer to an **mbuf** structure to a pointer to the head of an attached cluster.

The **MTOCL** macro can be used as in the following example:

```
caddr_t attcls;  
struct mbuf *m;  
attcls = (caddr_t) MTOCL(m);
```

## Execution Environment

The **MTOCL** macro can be called from either the process or interrupt environment.

### Related reference:

“M\_HASCL Macro for mbuf Kernel Services” on page 371

### Related information:

I/O Kernel Services

## MTOD Macro for mbuf Kernel Services

### Purpose

Converts a pointer to an **mbuf** structure to a pointer to the data stored in that **mbuf** structure.

### Syntax

```
#include <sys/mbuf.h>
```

```
MTOD ( m, type );
```

### Parameters

Item	Description
<i>m</i>	Identifies the address of an <b>mbuf</b> structure.
<i>type</i>	Indicates the type to which the resulting pointer should be cast.

### Description

The **MTOD** macro converts a pointer to an **mbuf** structure into a pointer to the data stored in the **mbuf** structure. This macro can be used as in the following example:

```
char *bufp;  
bufp = MTOD(m, char *);
```

### Execution Environment

The **MTOD** macro can be called from either the process or interrupt environment.

#### Related reference:

“DTOM Macro for mbuf Kernel Services” on page 106

#### Related information:

I/O Kernel Services

## M\_XMEMD Macro for mbuf Kernel Services

### Purpose

Returns the address of an **mbuf** cross-memory descriptor.

### Syntax

```
#include <sys/mbuf.h>  
#include <sys/xmem.h>
```

```
struct mbuf * m;  
M_XMEMD (m);
```

### Parameter



Item	Description
<i>m</i>	Specifies the address of the <b>mbuf</b> structure in question.

## Description

The **M\_XMEMD** macro returns the address of an **mbuf** cross-memory descriptor.

## Execution Environment

The **M\_XMEMD** macro can be called from either the process or interrupt environment.

## Example

The **M\_XMEMD** macro can be used as in the following example:

```
struct mbuf    *m;
struct xmem    *xmemd;
xmemd = M_XMEMD(m);
```

### Related information:

I/O Kernel Services

## mycpu Kernel Service

### Purpose

Gets the bind ID of the processor we are running on.

### Syntax

```
#include <sys/processor.h>
```

```
cpu_t myc ()
```

### Description

The **mycpu** kernel service returns the bind ID of the processor we are currently running on.

### Execution Environment

The **mycpu** kernel services can be called from either the process or interrupt environment. This routine must be called disabled. Otherwise, the calling thread might be preempted and resume execution on a different processor resulting in a stale value being returned.

### Return Values

The **mycpu** kernel service returns the bind ID of the current processor.

### Related reference:

“bindprocessor Kernel Service” on page 28

## n

The following kernel services begin with the with the letter n.

## nameToXfid() Kernel Service

### Purpose

Obtains the **xfid** value and attributes for a specific file name.

## Syntax

```
#include <sys/xfops.h>
#include <sys/vattr.h>

int nameToXfid(char *pathname,
               struct xfid *xfp,
               struct vattr *vap,
               long flags);
```

## Description

A kernel extension might need to convert a path name to an `xfid_t` structure. The `nameToXfid()` kernel service returns the `xfid` value for a specific path name.

## Parameters

### **pathname**

Full path name of the file for which an `xfid` value is needed.

### **xfp**

Pointer to an `xfid_t` structure to hold the `xfid` value that is set by this routine.

### **vap**

Pointer to a `vattr` structure to be entered by this routine. No attributes are set if the pointer is null.

### **flags**

Operation modifiers. This parameter must be set to zero.

## Return values

0 Indicates success. The `xfid` value and the optional `vattr` structure are returned.

### **ENOENT**

Name not found.

### **EPERM**

No permission for lookup.

### **EINVAL**

Invalid parameter is specified.

## net\_attach Kernel Service

### Purpose

Opens a communications I/O device handler.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <aixif/net_if.h>
#include <sys/comio.h>

int net_attach (kopen_ext, device_req, netid, netfpp)
struct kopen_ext * kopen_ext;
struct device_req * device_req;
struct netid_list * netid;
struct file ** netfpp;
```

### Parameters

Item	Description
<i>kopen_ext</i>	Specifies the device handler kernel open extension.
<i>device_req</i>	Indicates the address of the device description structure.
<i>netid</i>	Indicates the address of the network ID list.
<i>netfpp</i>	Specifies the address of the variable that will hold the returned file pointer.

## Description

The **net\_attach** kernel service opens the device handler specified by the *device\_req* parameter and then starts all the network IDs listed in the address specified by the *netid* parameter. The **net\_attach** service then sleeps and waits for the asynchronous start completion notifications from the **net\_start\_done** kernel service.

## Execution Environment

The **net\_attach** kernel service can be called from the process environment only.

## Return Values

Upon success, a value of 0 is returned and a file pointer is stored in the address specified by the *netfpp* parameter. Upon failure, the **net\_attach** service returns either the error codes received from the **fp\_opendev** or **fp\_ioctl** kernel service, or the value **ETIMEDOUT**. The latter value is returned when an open operation times out.

### Related reference:

“net\_detach Kernel Service”

“net\_start Kernel Service” on page 382

“net\_start\_done Kernel Service” on page 383

### Related information:

Network Kernel Services

## net\_detach Kernel Service

### Purpose

Closes a communications I/O device handler.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <aixif/net_if.h>
```

```
int net_detach ( netfp)
struct file *netfp;
```

### Parameter

Item	Description
<i>netfp</i>	Points to an open file structure obtained from the <b>net_attach</b> kernel service.

## Description

The **net\_detach** kernel service closes the device handler associated with the file pointer specified by the *netfp* parameter.

## Execution Environment

The **net\_detach** kernel service can be called from the process environment only.

## Return Values

The **net\_detach** service returns the value it obtains from the **fp\_close** service.

### Related reference:

“fp\_close Kernel Service” on page 145

“net\_attach Kernel Service” on page 378

### Related information:

Network Kernel Services

## net\_error Kernel Service

### Purpose

Handles errors for communication network interface drivers.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <net/if.h>
#include <sys/comio.h>
```

```
net_error ( ifp, error_code, netfp)
struct ifnet *ifp;
int error_code;
struct file *netfp;
```

### Parameters

Item	Description
<i>error_code</i>	Specifies the error code listed in the <b>/usr/include/sys/comio.h</b> file.
<i>ifp</i>	Specifies the address of the <b>ifnet</b> structure for the device with an error.
<i>netfp</i>	Specifies the file pointer for the device with an error.

## Description

The **net\_error** kernel service provides generic error handling for communications network interface (**if**) drivers. Network interface (**if**) kernel extensions call this service to trace errors and, in some instances, perform error recovery.

Errors traced include those:

- Received from the communications adapter drivers.
- Occurring during input and output packet processing.

## Execution Environment

The `net_error` kernel service can be called from either the process or interrupt environment.

## Return Values

The `net_error` service has no return values.

### Related reference:

“`net_attach` Kernel Service” on page 378

“`net_detach` Kernel Service” on page 379

### Related information:

Network Kernel Services

## `net_sleep` Kernel Service

### Purpose

Sleeps on the specified wait channel.

### Syntax

```
#include <sys/types.h>
```

```
#include <sys/errno.h>
```

```
#include <sys/pri.h>
```

```
net_sleep ( chan, flags )
```

```
int chan;
```

```
int flags;
```

### Parameters

Item	Description
<i>chan</i>	Specifies the wait channel to sleep upon.
<i>flags</i>	Sleep flags described in the <code>sleep</code> kernel service.

### Description

The `net_sleep` kernel service puts the caller to sleep waiting on the specified wait channel. If the caller holds the network lock, the `net_sleep` kernel service releases the lock before sleeping and reacquires the lock when the caller is awakened.

## Execution Environment

The `net_sleep` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates that the sleeping process was not awakened by a signal.
1	Indicates that the sleeper was awakened by a signal.

#### Related reference:

“net\_wakeup Kernel Service” on page 383

“sleep Kernel Service” on page 476

#### Related information:

Network Kernel Services

## net\_start Kernel Service

### Purpose

Starts network IDs on a communications I/O device handler.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <aixif/net_if.h>
#include <sys/comio.h>
```

```
struct file *net_start ( netfp, netid)
struct file *netfp;
struct netid_list *netid;
```

### Parameters

Item	Description
<i>netfp</i>	Specifies the file pointer of the device handler.
<i>netid</i>	Specifies the address of the network ID list.

### Description

The **net\_start** kernel service starts all the network IDs listed in the list specified by the *netid* parameter. This service then waits for the asynchronous notification of completion of starts.

### Execution Environment

The **net\_start** kernel service can be called from the process environment only.

### Return Values

The **net\_start** service uses the return value returned from a call to the **fp\_ioctl** service requesting the **CIO\_START** operation.

Item	Description
ETIMEDOUT	Indicates that the start for at least one network ID timed out waiting for start-done notifications from the device handler.

#### Related reference:

“fp\_ioctl Kernel Service” on page 151

“net\_attach Kernel Service” on page 378

“net\_start\_done Kernel Service” on page 383

#### Related information:

## net\_start\_done Kernel Service

### Purpose

Starts the done notification handler for communications I/O device handlers.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <aixif/net_if.h>
#include <sys/comio.h>
```

```
void net_start_done ( netid, sbp)
struct netid_list *netid;
struct status_block *sbp;
```

### Parameters

Item	Description
<i>netid</i>	Specifies the address of the network ID list for the device being started.
<i>sbp</i>	Specifies the status block pointer returned from the device handler.

### Description

The **net\_start\_done** kernel service is used to mark the completion of a network ID start operation. When all the network IDs listed in the *netid* parameter have been started, the **net\_attach** kernel service returns to the caller. The **net\_start\_done** service should be called when a **CIO\_START\_DONE** status block is received from the device handler. If the status block indicates an error, the start process is immediately aborted.

### Execution Environment

The **net\_start\_done** kernel service can be called from either the process or interrupt environment.

### Return Values

The **net\_start\_done** service has no return values.

#### Related reference:

“net\_attach Kernel Service” on page 378

“net\_start Kernel Service” on page 382

#### Related information:

CIO\_START\_DONE subroutine

Network Kernel Services

## net\_wakeup Kernel Service

### Purpose

Wakes up all sleepers waiting on the specified wait channel.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
net_wakeup ( chan)
int chan;
```

## Parameter

Item	Description
<i>chan</i>	Specifies the wait channel.

## Description

The `net_wakeup` service wakes up all network processes sleeping on the specified wait channel.

## Execution Environment

The `net_wakeup` kernel service can be called from either the process or interrupt environment.

## Return Values

The `net_wakeup` service has no return values.

### Related reference:

“`net_sleep` Kernel Service” on page 381

### Related information:

Network Kernel Services

## net\_xmit Kernel Service

### Purpose

Transmits data using a communications device handler .

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <aixif/net_if.h>
```

```
int net_xmit (ifp, m, netfp, lngth, m_ext)
struct ifnet * ifp;
struct mbuf * m;
struct file * netfp;
int lngth;
struct mbuf * m_ext;
```

### Parameters

Item	Description
<i>ifp</i>	Indicates an address of the <code>ifnet</code> structure for this interface.
<i>m</i>	Specifies the address of an <code>mbuf</code> structure containing the data to transmit.
<i>netfp</i>	Indicates the open file pointer obtained from the <code>net_attach</code> kernel service.
<i>lngth</i>	Indicates the total length of the buffer being transmitted.
<i>m_ext</i>	Indicates the address of an <code>mbuf</code> structure containing a write extension.

## Description

The `net_xmit` kernel service builds a `uio` structure and then invokes the `fp_rwuio` service to transmit a packet. The `net_xmit_trace` kernel service is an alternative for network interfaces that choose not to use the `net_xmit` kernel service.



## Execution Environment

The `net_xmit` kernel service can be called from either the process or interrupt environment.

## Return Values

Item	Description
0	Indicates that the packet was transmitted successfully.
ENOBUFS	Indicates that buffer resources were not available.

The `net_xmit` kernel service returns a value from the `fp_rwuio` service when an error occurs during a call to that service.

### Related reference:

“`fp_rwuio` Kernel Service” on page 166

“`net_xmit_trace` Kernel Service”

### Related information:

Network Kernel Services

## `net_xmit_trace` Kernel Service

### Purpose

Traces transmit packets.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
int net_xmit_trace ( ifp, mbuf )
struct ifnet *ifp;
struct mbuf *mbuf;
```

### Parameters

Item	Description
<i>ifp</i>	Designates the address of the <code>ifnet</code> structure for this interface.
<i>mbuf</i>	Designates the address of the <code>mbuf</code> structure to be traced.

### Description

The `net_xmit_trace` kernel service traces the data pointed to by the `mbuf` parameter. This kernel service was added for those network interfaces that choose not to use the `net_xmit` kernel service to transmit packets. An application program (the `iptrace` command) reads the trace data and writes it to a file for the `ipreport` command to interpret.

## Execution Environment

The `net_xmit_trace` kernel service can be called from either the process or interrupt environment.

## Return Values

The `net_xmit_trace` kernel service has no return values.

### Related reference:

“`net_xmit` Kernel Service” on page 384

**Related information:**

ipreport subroutine  
iptrace subroutine  
Network Kernel Services

## NLuprintf Kernel Service

### Purpose

Submits a request to print an internationalized message to a process' controlling terminal.

### Syntax

```
#include <sys/uprintf.h>
int NLuprintf (Uprintf)
struct uprintf *Uprintf;
```

### Parameters

Item	Description
<i>Uprintf</i>	Points to a <b>uprintf</b> request structure.

### Description

The **NLuprintf** kernel service submits a internationalized kernel message request with the **uprintf** request structure specified by the *Uprintf* parameter as input. Once the request has been successfully submitted, the **uprintfd** daemon retrieves, converts, formats, and writes the message described by the **uprintf** request structure to a process' controlling terminal.

The caller must initialize the **uprintf** request structure before calling the **NLuprintf** kernel service. Fields in the **uprintf** request structure use several constants. The following constants are defined in the */usr/include/sys/uprintf.h* file:

- **UP\_MAXSTR**
- **UP\_MAXARGS**
- **UP\_MAXCAT**
- **UP\_MAXMSG**

The **uprintf** request structure consists of the following fields:

## Field

Uprintf->upf\_defmsg

## Description

Points to a default message format. The default message format is a character string that contains either or both of two types of objects:

- Plain characters, which are copied to the message output stream
- Conversion specifications, each of which causes zero or more items to be fetched from the *Uprintf->arg* value parameter array

Each conversion specification consists of a % (percent sign) followed by a character that indicates the type of conversion to be applied:

- %** Performs no conversion. Prints a % character.
- d, i** Accepts an integer value and converts it to signed decimal notation.
- u** Accepts an integer value and converts it to unsigned decimal notation.
- o** Accepts an integer value and converts it to unsigned octal notation.
- x** Accepts an integer value and converts it to unsigned hexadecimal notation.
- c** Accepts and prints a **char** value.
- s** Accepts a value as a string (character pointer). Characters from the string are printed until a \0 (null character) is encountered.

Field-width or precision conversion specifications are not supported.

The maximum length of the default message-format string pointed to by the Uprintf->upf\_defmsg field is the number of characters specified by the UP\_MAXSTR constant. The Uprintf->upf\_defmsg field must be a nonnull character.

The default message format is used in constructing the kernel message if the message format described by the Uprintf->upf\_NLsetno and Uprintf->upf\_NLmsgno fields cannot be retrieved from the message catalog specified by Uprintf->upf\_NLcatname. The conversion specifications contained within the default message format should match those contained in the message format specified by the upf\_NLsetno and upf\_NLmsgno fields.

Uprintf->upf\_arg[UP\_MAXARGS]

Specifies from zero to the number of value parameters specified by the UP\_MAXARGS constant. A *Value* parameter may be an integer value, a character value, or a string value (character pointer). Strings are limited in length to the number of characters specified by the UP\_MAXSTR constant. String value parameters must be nonnull characters. The number, type, and order of items in the *Value* parameter array should match the conversion specifications within the message format string.

Uprintf->upf\_NLcatname

Points to the message catalog file name. If the catalog file name referred to by the Uprintf->upf\_NLcatname field begins with a / (slash), it is assumed to be an absolute path name. If the catalog file name is not an absolute path name, the process environment determines the directory paths to search. The maximum length of the catalog file name is limited to the number of characters specified by the UP\_MAXCAT constant. The value of the Uprintf->upf\_NLcatname field must be a nonnull character.

Uprintf->upf\_NLsetno

Specifies the set ID.

Field	Description
Uprintf->upf_NLmsgno	Specifies the message ID. The Uprintf->upf_NLsetno and Uprintf->upf_NLmsgno fields specify a particular message format string to be retrieved from the message catalog specified by the Uprintf->upf_NLcatname field.  The maximum length of the constructed kernel message is limited to the number of characters specified by the <b>UP_MAXMSG</b> constant. Messages larger than the number of characters specified by the <b>UP_MAXMSG</b> constant are discarded.

## Execution Environment

The **NLuprntf** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
ENOMEM	Indicates that memory is not available to buffer the request.
ENODEV	Indicates that a controlling terminal does not exist for the process.
ESRCH	Indicates the <b>uprntfd</b> daemon is not active. No requests may be submitted.
EINVAL	Indicates that the message catalog file-name pointer is null or the catalog file name is greater than the number of characters specified by the <b>UP_MAXCAT</b> constant.
EINVAL	Indicates that a string-value parameter pointer is null or the string-value parameter is greater than the number of characters specified by the <b>UP_MAXCAT</b> constant.
EINVAL	Indicates one of the following: <ul style="list-style-type: none"> <li>• Default message format pointer is null.</li> <li>• Number of characters in the default message format is greater than the number specified by the <b>UP_MAXSTR</b> constant.</li> <li>• Number of conversion specifications contained within the default message format is greater than the number specified by the <b>UP_MAXARGS</b> constant.</li> </ul>

### Related reference:

“uprntf Kernel Service” on page 528

### Related information:

uprntfd subroutine

Process and Exception Management Kernel Services

## ns\_add\_demux Network Kernel Service

### Purpose

Adds a demuxer for the specified type of network interface.

### Syntax

```
#include <sys/ndd.h>
#include <sys/cdli.h>
```

```
int ns_add_demux (ndd_type, demux)
    u_long ndd_type;
    struct ns_demuxer * demux;
```

### Parameters

Item	Description
<i>ndd_type</i>	Specifies the interface type of the demuxer to be added.
<i>demux</i>	Specifies the pointer to an <b>ns_demux</b> structure that defines the demuxer.

## Description

The **ns\_add\_demux** network service adds the specified demuxer to the list of available network demuxers. Only one demuxer per network interface type can exist. An interface type describes a certain class of network devices that have the same characteristics (such as ethernet or token ring). The values of the *ndd\_type* parameter listed in the **/usr/include/sys/ndd.h** file are the numbers defined by Simple Network Management Protocol (SNMP). If the desired type is not in the **ndd.h** file, the SNMP value should be used if it is defined. Otherwise, any undefined type above **NDD\_MAX\_TYPE** may be used.

**Note:** The **ns\_demux** structure must be allocated and pinned by the network demuxer.

## Examples

The following example illustrates the **ns\_add\_demux** network service:

```
struct ns_demux demuxer;
bzero (&demuxer, sizeof (demuxer));
demuxer.nd_add_filter = eth_add_filter;
demuxer.nd_del_filter = eth_del_filter;
demuxer.nd_add_status = eth_add_status;
demuxer.nd_del_status = eth_del_status;
demuxer.nd_receive = eth_receive;
demuxer.nd_status = eth_status;
demuxer.nd_response = eth_response;
demuxer.nd_use_nsdnx = 1;
ns_add_demux(NDD_IS088023, &demuxer);
```

## Return Values

Item	Description
0	Indicates the operation was successful.
EEXIST	Indicates a demuxer already exists for the given type.

## Related reference:

“ns\_del\_demux Network Service” on page 395

## ns\_add\_filter Network Service

### Purpose

Registers a receive filter to enable the reception of packets.

### Syntax

```
#include <sys/cdli.h>
#include <sys/ndd.h>
```

```
int ns_add_filter (nddp, filter, len, ns_user)
    struct ndd * nddp;
    caddr_t filter;
    int len;
    struct ns_user * ns_user;
```

### Parameters

Item	Description
<i>nddp</i>	Specifies the <b>ndd</b> structure to which this add request applies.
<i>filter</i>	Specifies the pointer to the receive filter.
<i>len</i>	Specifies the length in bytes of the receive filter to which the <i>filter</i> parameter points.
<i>ns_user</i>	Specifies the pointer to a <b>ns_user</b> structure that defines the user.

## Description

The **ns\_add\_filter** network service registers a receive filter for the reception of packets and enables a network demuxer to route packets to the appropriate users. The **add** request is passed on to the **nd\_add\_filter** function of the demuxer for the specified NDD. The caller of the **ns\_add\_filter** network service is responsible for relinquishing filters before calling the **ns\_free** network service.

## Examples

The following example illustrates the **ns\_add\_filter** network service:

```
struct ns_8022 d1;
struct ns_user ns_user;

d1.filtertype = NS_LLC_DSAP_SNAP;
d1.dsap = 0xaa;
d1.orgcode[0] = 0x0;
d1.orgcode[1] = 0x0;
d1.orgcode[2] = 0x0;
d1.ethertype = 0x0800;

ns_user.isr = NULL;
ns_user.isr_data = NULL;
ns_user.protoq = &ipintrq;
ns_user.netisr = NETISR_IP;
ns_user.ifp = ifp;
ns_user.pkt_format = NS_PROTO_SNAP;

ns_add_filter(nddp, &d1, sizeof(d1), &ns_user);
```

There are two ways a user (that is, the entity that is interested in receiving incoming packets) can be invoked when a packet arrives. In the first method, a protocol queue can be defined in which incoming packets are queued upon receipt, and the specified *netisr* is scheduled to let the user know that there are new packets in the queue. For example, the preceding code assumes a network interrupt service request (*netisr*) with the name **NETISR\_IP** has been defined. When a packet arrives for the specified user, the packet is queued on the specified protocol queue (in this case, **ipintrq**) and the **NETISR\_IP** request is scheduled to be executed. Because of its complexity, this mode is not currently being used by any network user.

The preferred way of receiving incoming packets is by registering an interrupt service request (*isr*) function that handles incoming packets; **ns\_user.isr** points to the function that will get invoked whenever a packet that matches the specified filter arrives. This function should expect the following four arguments:

```
void isr (ndd_t *nddp, mbuf *m, caddr_t macp, caddr_t extp)
```

where

Item	Description
<i>nddp</i>	Pointer to the <b>ndd</b> structure representing the adapter where the packet was received.
<i>m</i>	Pointer to the <b>mbuf</b> structure representing the packet that was received.
<i>macp</i>	Pointer to the start of the MAC header of the packet that was received.
<i>extp</i>	Pointer to the (optional) structure specified in <b>ns_user.isr_data</b> , or NULL if none was specified.

In the following code, the function **bpf\_cdli\_tap** will be called when a new packet arrives; a pointer to the **bp** structure will be passed as the fourth parameter when **bpf\_cdli\_tap** is called.

```
dl.filtertype = NS_TAP;

ns_user.isr = bpf_cdli_tap;
ns_user.isr_data = (caddr_t) bp;
ns_user.protoq = (struct ifqueue *) NULL;
ns_user.netisr = 0;
ns_user.ifp = (struct ifnet *) NULL;
ns_user.pkt_format = NS_INCLUDE_MAC;
```

**Note:** Both modes of receiving packets are mutually exclusive. In other words, if the **ns\_user.protoq** member is non-null, the protocol queue method is used; otherwise, the direct isr function method is used, and the **ns\_user.isr** function pointer must be a valid function pointer.

In both cases, **ns\_user.ifp** can optionally point to the **ifnet** structure of the interface where the packets will be received. If it is non-null, the state of the interface will be verified when a packet is received. If the interface is not up, the packet will be dropped and it will not be delivered to the user. If the interface is up, the statistics for the number of received packets will be incremented, and the ifp will be saved in the packet's **mbuf** structure's **m\_pkthdr.rcvif** field.

The **ns\_user.pkt\_format** member determines how much of the MAC header the user is interested in receiving. Its possible values are:

Item	Description
<b>NS_PROTO</b>	Do not include the LLC header (but include the SNAP header, if there is one).
<b>NS_PROTO_SNAP</b>	Do not include the LLC SNAP header (that is, remove the entire MAC header and deliver only the data).
<b>NS_INCLUDE_LLC</b>	Include the LLC header.
<b>NS_INCLUDE_MAC</b>	Include the entire MAC header.

Item	Description
NS_HANDLE_HEADERS	Instead of passing the specified <code>ns_user.isr_data</code> structure by itself, build an <code>isr_data_ext</code> structure containing header information, as well as a pointer to the specified <code>ns_user.isr_data</code> . These are the fields that will be set in the <code>isr_data_ext</code> structure:
<code>isr_data_ext.isr_data</code>	Pointer to the structure passed as <code>ns_user.isr_data</code> .
<code>isr_data_ext.dstp</code>	Pointer to the destination MAC address.
<code>isr_data_ext.dstlen</code>	Length of the destination MAC address.
<code>isr_data_ext.srcp</code>	Pointer to the source MAC address.
<code>isr_data_ext.seclen</code>	Length of the source MAC address.
<code>isr_data_ext.segp</code>	Pointer to the routing segment.
<code>isr_data_ext.seglen</code>	Length of the routing segment.
<code>isr_data_ext.llcp</code>	Pointer to the LLC.
<code>isr_data_ext.llclen</code>	Length of the LLC.
	It is possible to combine <code>NS_HANDLE_HEADERS</code> with one of the other flags by means of a logical OR operator (for example, <code>ns_user.pkt_format = NS_INCLUDE_MAC   NS_HANDLE_HEADERS</code> ). The other flags, however, are mutually exclusive.

## Return Values

Item	Description
0	Indicates the operation was successful.

The network demuxer may supply other return values.

### Related reference:

“ns\_del\_filter Network Service” on page 396

## ns\_add\_status Network Service

### Purpose

Adds a status filter for the routing of asynchronous status.

### Syntax

```
#include <sys/cdli.h>
#include <sys/ndd.h>
```

```
int ns_add_status (nddp, statfilter, len, ns_statuser)
    struct ndd * nddp;
    caddr_t statfilter;
    int len;
    struct ns_statuser * ns_statuser;
```

### Parameters



Item	Description
<i>nddp</i>	Specifies a pointer to the <b>ndd</b> structure to which this add request applies.
<i>statfilter</i>	Specifies a pointer to the status filter.
<i>len</i>	Specifies the length, in bytes, of the value of the <i>statfilter</i> parameter.
<i>ns_statuser</i>	Specifies a pointer to an <b>ns_statuser</b> structure that defines this user.

## Description

The **ns\_add\_status** network service registers a status filter. The add request is passed on to the **nd\_add\_status** function of the demuxer for the specified network device driver (NDD). This network service enables the user to receive asynchronous status information from the specified device.

**Note:** The user's status processing function is specified by the *isr* field of the **ns\_statuser** structure. The network demuxer calls the user's status processing function directly when asynchronous status information becomes available. Consequently; the status processing function cannot be a scheduled routine. The caller of the **ns\_add\_status** network service is responsible for relinquishing status filters before calling the **ns\_free** network service.

## Examples

The following example illustrates the **ns\_add\_status** network service:

```
struct ns_statuser  user;
struct ns_com_status  filter;

filter.filtertype = NS_STATUS_MASK;
filter.mask = NDD_HARD_FAIL;
filter.sid = 0;
user.isr = status_fn;
user.isr_data = whatever_makes_sense;

error = ns_add_status(nddp, &filter, sizeof(filter), &user);
```

## Return Values

Item	Description
0	Indicates the operation was successful.

The network demuxer may supply other return values.

### Related reference:

“ns\_del\_status Network Service” on page 397

## ns\_alloc Network Service

### Purpose

Allocates use of a network device driver (NDD).

### Syntax

```
#include <sys/ndd.h>
```

```
int ns_alloc (nddname, nddpp)
    char * nddname;
    struct ndd ** nddpp;
```

### Parameters

Item	Description
<i>nddname</i>	Specifies the device name to be allocated.
<i>nddpp</i>	Indicates the address of the pointer to a <b>ndd</b> structure.

## Description

The **ns\_alloc** network service searches the Network Service (NS) device chain to find the device driver with the specified *nddname* parameter. If the service finds a match, it increments the reference count for the specified device driver. If the reference count is incremented to 1, the **ndd\_open** subroutine specified in the **ndd** structure is called to open the device driver.

## Examples

The following example illustrates the **ns\_alloc** network service:

```
struct ndd *nddp;
error = ns_alloc("en0", &nddp);
```

## Return Values

If a match is found and the **ndd\_open** subroutine to the device is successful, a pointer to the **ndd** structure for the specified device is stored in the *nddpp* parameter. If no match is found or the open of the device is unsuccessful, a non-zero value is returned.

Item	Description
0	Indicates the operation was successful.
ENODEV	Indicates an invalid network device.
ENOENT	Indicates no network demuxer is available for this device.

The **ndd\_open** routine may specify other return values.

### Related reference:

“ns\_free Network Service” on page 398

## ns\_attach Network Service

### Purpose

Attaches a network device to the network subsystem.

### Syntax

```
#include <sys/ndd.h>
```

```
int ns_attach (nddp)
    struct ndd * nddp;
```

### Parameters

Item	Description
<i>nndp</i>	Specifies a pointer to an <b>nnd</b> structure describing the device to be attached.

## Description

The **ns\_attach** network service places the device into the available network service (NS) device chain. The network device driver (NDD) should be prepared to be opened after the **ns\_attach** network service is called.

**Note:** The **nnd** structure is allocated and initialized by the device. It should be pinned.

## Examples

The following example illustrates the **ns\_attach** network service:

```
struct nnd nnd;
nnd.nnd_name = "en0";
nnd.nnd_addrLen = 6;
nnd.nnd_hdrLen = 14;
nnd.nnd_mtu = ETHERMTU;
nnd.nnd_mintu = 60;
nnd.nnd_type = NDD_ETHER;
nnd.nnd_flags =
    NDD_BROADCAST | NDD_SIMPLEX;
nnd.nnd_open = entopen;
nnd.nnd_output = entwrite;
nnd.nnd_ctl = entctl;
nnd.nnd_close = entclose;
.
.
.
ns_attach(&nnd);
```

## Return Values

Item	Description
0	Indicates the operation was successful.
EEXIST	Indicates the device is already in the available NS device chain.

### Related reference:

“ns\_detach Network Service” on page 398

## ns\_del\_demux Network Service

### Purpose

Deletes a demuxer for the specified type of network interface.

### Syntax

```
#include <sys/nnd.h>
```

```
int ns_del_demux (nnd_type)
    u_long nnd_type;
```

### Parameters

Item	Description
<i>ndd_type</i>	Specifies the network interface type of the demuxer that is to be deleted.

## Description

If the demuxer is not currently in use, the **ns\_del\_demux** network service deletes the specified demuxer from the list of available network demuxers. A demuxer is in use if a network device driver (NDD) is open for the demuxer.

## Examples

The following example illustrates the **ns\_del\_demux** network service:

```
ns_del_demux(NDD_IS088023);
```

## Return Values

Item	Description
0	Indicates the operation was successful.
ENOENT	Indicates the demuxer of the specified type does not exist.

## Related reference:

“ns\_add\_demux Network Kernel Service” on page 388

## ns\_del\_filter Network Service

### Purpose

Deletes a receive filter.

### Syntax

```
#include <sys/cdli.h>
#include <sys/ndd.h>
```

```
int ns_del_filter (nddp, filter, len)
    struct ndd * nddp;
    caddr_t filter;
    int len;
```

### Parameters

Item	Description
<i>nddp</i>	Specifies the <b>ndd</b> structure that this delete request is for.
<i>filter</i>	Specifies the pointer to the receive filter.
<i>len</i>	Specifies the length in bytes of the receive filter.

## Description

The **ns\_del\_filter** network service deletes the receive filter from the corresponding network demuxer. This disables packet reception for packets that match the filter. The delete request is passed on to the **nd\_del\_filter** function of the demuxer for the specified network device driver (NDD).

## Examples

The following example illustrates the **ns\_del\_filter** network service:

```

struct ns_8022 dl;

dl.filtertype = NS_LLC_DSAP_SNAP;
dl.dsap = 0xaa;
dl.orgcode[0] = 0x0;
dl.orgcode[1] = 0x0;
dl.orgcode[2] = 0x0;
dl.ethertype = 0x0800;
ns_del_filter(nddp, &dl, sizeof(dl));

```

## Return Values

Item	Description
0	Indicates the operation was successful.

The network demuxer may supply other return values.

### Related reference:

“ns\_add\_filter Network Service” on page 389

“ns\_alloc Network Service” on page 393

## ns\_del\_status Network Service Purpose

Deletes a previously added status filter.

### Syntax

```

#include <sys/cdli.h>
#include <sys/ndd.h>

```

```

int ns_del_status (nddp, statfilter, len)
    struct ndd * nddp;
    caddr_t statfilter;
    int len;

```

### Parameters

Item	Description
<i>nddp</i>	Specifies the pointer to the <b>ndd</b> structure to which this delete request applies.
<i>statfilter</i>	Specifies the pointer to the status filter.
<i>len</i>	Specifies the length, in bytes, of the value of the <i>statfilter</i> parameter.

### Description

The **ns\_del\_status** network service deletes a previously added status filter from the corresponding network demuxer. The delete request is passed on to the **nd\_del\_status** function of the demuxer for the specified network device driver (NDD). This network service disables asynchronous status notification from the specified device.

### Examples

The following example illustrates the **ns\_del\_status** network service:

```

error = ns_add_status(nddp, &filter,
sizeof(filter));

```

## Return Values

Item	Description
0	Indicates the operation was successful.

The network demuxer may supply other return values.

### Related reference:

“ns\_add\_status Network Service” on page 392

## ns\_detach Network Service

### Purpose

Removes a network device from the network subsystem.

### Syntax

```
#include <sys/ndd.h>
```

```
int ns_detach (nddp)
    struct ndd * nddp;
```

### Parameters

Item	Description
<i>nddp</i>	Specifies a pointer to an <b>ndd</b> structure describing the device to be detached.

### Description

The **ns\_detach** service removes the **ndd** structure from the chain of available NS devices.

### Examples

The following example illustrates the **ns\_detach** network service:

```
ns_detach(nddp);
```

## Return Values

Item	Description
0	Indicates the operation was successful.
ENOENT	Indicates the specified <i>ndd</i> structure was not found.
EBUSY	Indicates the network device driver (NDD) is currently in use.

### Related reference:

“ns\_attach Network Service” on page 394

## ns\_free Network Service

### Purpose

Relinquishes access to a network device.

### Syntax

```
#include <sys/ndd.h>
```

```
void ns_free (nddp)
    struct ndd * nddp;
```

## Parameters

Item	Description
<i>nndp</i>	Specifies the <b>nnd</b> structure of the network device that is to be freed from use.

## Description

The **ns\_free** network service relinquishes access to a network device. The **ns\_free** network service also decrements the reference count for the specified **nnd** structure. If the reference count becomes 0, the **ns\_free** network service calls the **nnd\_close** subroutine specified in the **nnd** structure.

## Examples

The following example illustrates the **ns\_free** network service:

```
struct nnd *nndp;
ns_free(nndp);
```

## Files

Item	Description
<code>net/cdli.c</code>	

### Related reference:

“**ns\_alloc** Network Service” on page 393

## p

The following kernel services begin with the with the letter p.

### **\_\_pag\_getid** System Call

#### Purpose

Invokes the **kcred\_getpagid** kernel service and returns the PAG identifier for that PAG name.

#### Syntax

```
int __pag_getid (name)
char *name;
```

#### Description

Given a PAG type name, the **\_\_pag\_getid** invokes the **kcred\_getpagid** kernel service and returns the PAG identifier for that PAG name.

#### Parameters

Item	Description
<i>name</i>	A <b>char *</b> value which references a NULL-terminated string of not more than <b>PAG_NAME_LENGTH_MAX</b> characters.

#### Return Values

If successful, a value greater than or equal to 0 is returned and represents the PAG type. This value may be used in subsequent calls to other PAG system calls that require a *type* parameter on input. If unsuccessful, -1 is returned and the **errno** global variable is set to a value reflecting the cause of the error.

#### Error Codes

Item	Description
ENOENT	The <i>name</i> parameter doesn't refer to an existing PAG type.
ENAMETOOLONG	The <i>name</i> parameter refers to a string that is longer than PAG_NAME_LENGTH_MAX.

#### Related reference:

“\_\_pag\_getname System Call”

“\_\_pag\_setname System Call” on page 401

“kcred\_getpagid Kernel Service” on page 246

## \_\_pag\_getname System Call

### Purpose

Retrieves the name of a PAG type.

### Syntax

```
int __pag_getname (type, buf, size)
int type;
char *buf;
int size;
```

### Description

The `__pag_getname` system call retrieves the name of a PAG type given its integer value by invoking the `kcred_getpagname` kernel service with the given parameters.

### Parameters

Item	Description
<i>type</i>	A numerical PAG identifier.
<i>buf</i>	A <code>char *</code> value that points to an array at least PAG_NAME_LENGTH_MAX+1 bytes in length.
<i>size</i>	An <code>int</code> value that gives the size of <i>buf</i> in bytes.

### Return Values

If successful, 0 is returned and the *buf* parameter contains the PAG name associated with the *type* parameter. If unsuccessful, -1 is returned and the `errno` global variable is set to a value reflecting the cause of the error.

### Error Codes

Item	Description
EINVAL	The value of the <i>type</i> parameter is less than 0 or greater than the maximum PAG identifier.
ENOENT	There is no PAG associated with the <i>type</i> parameter.
ENOSPC	The value of the <i>size</i> parameter is insufficient to hold the PAG name and its terminating NULL character.

#### Related reference:

“\_\_pag\_getvalue System Call”

“\_\_pag\_setname System Call” on page 401

“kcred\_getpagname Kernel Service” on page 247

## \_\_pag\_getvalue System Call

### Purpose

Invokes the `kcred_getpag` kernel service and returns the PAG value.



## Syntax

```
int __pag_getvalue (type)
int type;
```

## Description

Given a PAG type, the `__pag_getvalue` system call invokes the `kcred_getpag` kernel service and returns the PAG value for the value of the `type` parameter.

## Parameters

Item	Description
<code>type</code>	An <code>int</code> value indicating the desired PAG.

## Return Values

If successful, the value of the PAG (or 0 when there is no value for that PAG type) is returned. If unsuccessful, -1 is returned and the `errno` global variable is set to a value reflecting the cause of the error.

## Error Codes

Item	Description
EINVAL	The <code>type</code> parameter is less than 0 or greater than the maximum PAG type value.
ENOENT	The <code>type</code> parameter doesn't reference an existing PAG type.

**Note:** It is not an error for a defined PAG to not have a value in the current process' credentials.

### Related reference:

- “`__pag_getid` System Call” on page 399
- “`__pag_setvalue` System Call” on page 402
- “`kcred_getpagname` Kernel Service” on page 247

## `__pag_setname` System Call

### Purpose

Invokes the `kcred_setpagname` kernel service and returns the PAG type identifier.

## Syntax

```
int __pag_setname (name, flags)
char *name;
int flags;
```

## Description

The `__pag_setname` system call invokes the `kcred_setpagname` kernel service to register the name of a PAG and returns the PAG type identifier. The value of the `func` parameter to `kcred_setpagname` will be NULL. The other parameters to this system call are the same as with the underlying kernel service. This system call requires the `SYS_CONFIG` privilege.

## Parameters

Item	Description
<i>name</i>	A <b>char</b> * value giving the symbolic name of the requested PAG.
<i>flags</i>	Either PAG_UNIQUEVALUE or PAG_MULTIVALUED 1 .

## Return Values

A return value greater than or equal to 0 is the PAG type associated with the *name* parameter. This value may be used with other PAG-related system calls which require a numerical PAG identifier. If unsuccessful, -1 is returned and the **errno** global variable is set to indicate the cause of the error.

## Error Codes

Item	Description
ENOSPC	The PAG name table is full.
EEXIST	The named PAG type already exists in the table, and the <i>flags</i> and <i>func</i> parameters do not match their previous values.
EPERM	The calling process does not have the SYS_CONFIG privilege.

### Related reference:

“\_\_pag\_getname System Call” on page 400

“\_\_pag\_setvalue System Call”

“kcred\_setpagname Kernel Service” on page 252

## \_\_pag\_setvalue System Call

### Purpose

Invokes the **kcred\_setpag** kernel service and sets the value of PAG type to *pag*.

### Syntax

```
int __pag_setvalue (type, pag)
int type;
int pag;
```

### Description

Given a PAG type and value, the **\_\_pag\_setvalue** system call invokes the **kcred\_setpag** kernel service and sets the value of PAG type to *pag*. This system call requires the SET\_PROC\_DAC privilege.

### Parameters

Item	Description
<i>type</i>	An <b>int</b> value indicating the desired PAG.
<i>pag</i>	An <b>int</b> value containing the new PAG value.

## Return Values

If successful, 0 is returned. If unsuccessful, -1 is returned and the **errno** global variable is set to a value reflecting the cause of the error.

## Error Codes

Item	Description
ENOENT	The <i>type</i> parameter doesn't reference an existing PAG type.
EINVAL	The value of <i>pag</i> is -1.
EPERM	The calling process lacks the appropriate privilege.

**Related reference:**

- “\_\_pag\_getvalue System Call” on page 400
- “\_\_pag\_setname System Call” on page 401
- “kcred\_setpagname Kernel Service” on page 252

## panic Kernel Service

### Purpose

Crashes the system.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
panic ( s )
char *s;
```

### Parameter

Item	Description
s	Points to a character string to be written to the error log.

### Description

The **panic** kernel service is called when a catastrophic error occurs and the system can no longer continue to operate. The **panic** service performs these two actions:

- Writes the character string pointed to by the *s* parameter to the error log.
- Performs a system dump.

The system halts after the dump. You should wait for the dump to complete, reboot the system, and then save and analyze the dump.

### Execution Environment

The **panic** kernel service can be called from either the process or interrupt environment.

### Return Values

The **panic** kernel service has no return values.

**Related information:**

RAS Kernel Services

## pci\_cfgw Kernel Service

### Purpose

Reads and writes PCI bus slot configuration registers.

### Syntax

```
#include <sys/mdio.h>
```

```
int pci_cfgw(bid, md, write_flag)
int bid;
struct mdio *md;
int write_flag;
```

## Description

The `pci_cfgw` kernel service provides serialized access to the configuration registers for a PCI bus. To ensure data integrity in a multi-processor environment, a lock is required before accessing the configuration registers. Depending on the value of the `write_flag` parameter, a read or write to the configuration register is performed at offset `md_addr` for the device identified by `md_sla`.

The `pci_cfgw` kernel service provides for kernel extensions the same services as the `MIOPCFGET` and `MIOPCFPUT` ioctls provides for applications. The `pci_cfgw` kernel service can be called from either the process or the interrupt environment.

## Parameters

Item	Description
<i>bid</i>	Specifies the bus identifier.
<i>md</i>	Specifies the address of the <i>mdio</i> structure. The <i>mdio</i> structure contains the following fields: <ul style="list-style-type: none"> <li><i>md_addr</i> Starting offset of the configuration register to access (0 to 0xFF for PCI/PCI-X, and 0 to 0xFFF for PCI-E).</li> <li><i>md_data</i> Pointer to the data buffer.</li> <li><i>md_size</i> Number of items of size specified by the <i>md_incr</i> parameter. The maximum size is 256 bytes for PCI/PCI-X, and 4096 for PCI-E.</li> <li><i>md_incr</i> Access types, <code>MV_BYTE</code>, <code>MV_WORD</code>, or <code>MV_SHORT</code>.</li> <li><i>md_sla</i> Device Number and Function Number. (Device Number * 8) + Function.</li> </ul>
<i>write_flag</i>	Set to 1 for write and 0 for read.

## Return Values

Returns 0 for successful completion.

Item	Description
ENOMEM	Indicates no memory could be allocated.
EINVAL	Indicated that the bus, device/function, or size is not valid.
EPERM	Indicates that the platform does not allow the requested operation

### Related information:

Machine Device Driver

## pfctlinput Kernel Service Purpose

Invokes the `ctlinput` function for each configured protocol.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/domain.h>
```

```
void pfctlinput ( cmd, sa)
int cmd;
struct sockaddr *sa;
```

## Parameters

Item	Description
<i>cmd</i>	Specifies the command to pass on to protocols.
<i>sa</i>	Indicates the address of a <code>sockaddr</code> structure that is passed to the protocols.

## Description

The `pfctlinput` kernel service searches through the protocol switch table of each configured domain and invokes the protocol `ctlinput` function if defined. Both the *cmd* and *sa* parameters are passed as parameters to the protocol function.

## Execution Environment

The `pfctlinput` kernel service can be called from either the process or interrupt environment.

## Return Values

The `pfctlinput` service has no return values.

### Related information:

Network Kernel Services

Understanding Socket Header Files

## pffindproto Kernel Service Purpose

Returns the address of a protocol switch table entry.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/domain.h>
```

```
struct protosw *pffindproto (family, protocol, type)
int family;
int protocol;
int type;
```

## Parameters

Item	Description
<i>family</i>	Specifies the address family for which to search.
<i>protocol</i>	Indicates the protocol within the address family.
<i>type</i>	Specifies the type of socket (for example, <code>SOCK_RAW</code> ).

## Description

The `pffindproto` kernel service first searches the domain switch table for the address family specified by the *family* parameter. If found, the `pffindproto` service then searches the protocol switch table for that domain and checks for matches with the *type* and *protocol* parameters.

If a match is found, the **pfindproto** service returns the address of the protocol switch table entry. If the *type* parameter is set to **SOCK\_RAW**, the **pfindproto** service returns the first entry it finds with protocol equal to 0 and type equal to **SOCK\_RAW**.

## Execution Environment

The **pfindproto** kernel service can be called from either the process or interrupt environment.

## Return Values

The **pfindproto** service returns a null value if a protocol switch table entry was not found for the given search criteria. Upon success, the **pfindproto** service returns the address of a protocol switch table entry.

### Related information:

Network Kernel Services

Understanding Socket Header Files

## pgsignal Kernel Service

### Purpose

Sends a signal to all of the processes in a process group.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
void pgsignal ( pid, sig)
pid_t pid;
int sig;
```

### Parameters

Item	Description
------	-------------

<i>pid</i>	Specifies the process ID of a process in the group of processes to receive the signal.
<i>sig</i>	Specifies the signal to send.

### Description

The **pgsignal** kernel service sends a signal to each member in the process group to which the process identified by the *pid* parameter belongs. The *pid* parameter must be the process identifier of the member of the process group to be sent the signal. The *sig* parameter specifies which signal to send.

Device drivers can get the value for the *pid* parameter by using the **getpid** kernel service. This value is the process identifier for the currently executing process.

The **sigaction** subroutine contains a list of the valid signals.

## Execution Environment

The **pgsignal** kernel service can be called from either the process or interrupt environment.

## Return Values

The **pgsignal** service has no return values.

### Related reference:

“getpid Kernel Service” on page 190

“pidsig Kernel Service”

**Related information:**

sigaction subroutine

Process and Exception Management Kernel Services

## **pidsig Kernel Service**

### **Purpose**

Sends a signal to a process.

### **Syntax**

```
#include <sys/types.h>
```

```
#include <sys/errno.h>
```

```
void pidsig ( pid, sig)
```

```
pid_t pid;
```

```
int sig;
```

### **Parameters**

Item	Description
------	-------------

<i>pid</i>	Specifies the process ID of the receiving process.
------------	--

<i>sig</i>	Specifies the signal to send.
------------	-------------------------------

### **Description**

The **pidsig** kernel service sends a signal to a process. The *pid* parameter must be the process identifier of the process to be sent the signal. The *sig* parameter specifies the signal to send. See the **sigaction** subroutine for a list of the valid signals.

Device drivers can get the value for the *pid* parameter by using the **getpid** kernel service. This value is the process identifier for the currently executing process.

The **pidsig** kernel service can be called from an interrupt handler execution environment if the process ID is known.

### **Execution Environment**

The **pidsig** kernel service can be called from either the process or interrupt environment.

### **Return Values**

The **pidsig** service has no return values.

**Related reference:**

“getpid Kernel Service” on page 190

“pgsignal Kernel Service” on page 406

**Related information:**

sigaction subroutine

Process and Exception Management Kernel Services

## pin Kernel Service

### Purpose

Pins the address range in the system (kernel) space.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/pin.h>
```

```
int pin ( addr, length)
caddr_t addr;
int length;
```

### Parameters

Item	Description
<i>addr</i>	Specifies the address of the first byte to pin.
<i>length</i>	Specifies the number of bytes to pin.

### Description

The **pin** service pins the real memory pages touched by the address range specified by the *addr* and *length* parameters in the system (kernel) address space. It pins the real-memory pages to ensure that page faults do not occur for memory references in this address range. The **pin** service increments the pin count for each real-memory page. While the pin count is nonzero, the page cannot be paged out of real memory.

The **pin** routine pins either the entire address range or none of it. Only a limited number of pages can be pinned in the system. If there are not enough unpinned pages in the system, the **pin** service returns an error code.

**Note:** If the requested range is not aligned on a page boundary, then memory outside this range is also pinned. This is because the operating system pins only whole pages at a time.

The **pin** service can only be called for addresses within the system (kernel) address space. The **xmempin** service should be used for addresses within kernel or user space.

### Execution Environment

The **pin** kernel service can be called from the process environment only.

### Return Values

Item	Description
0	Indicates successful completion.
EINVAL	Indicates that the value of the <i>length</i> parameter is negative or 0. Otherwise, the area of memory beginning at the address of the first byte to pin (the <i>addr</i> parameter) and extending for the number of bytes specified by the <i>length</i> parameter is not defined.
EIO	Indicates that a permanent I/O error occurred while referencing data.
ENOMEM	Indicates that the <b>pin</b> service was unable to pin due to insufficient real memory or exceeding the systemwide pin count.
ENOSPC	Indicates insufficient file system or paging space.

### Related reference:

“xmempin Kernel Service” on page 605



“xmemunpin Kernel Service” on page 606

#### Related information:

Understanding Execution Environments

Memory Kernel Services

## pin\_context\_stack or unpin\_context\_stack Kernel Service

### Purpose

Pins and unpins hidden kernel stack region.

### Syntax

```
#include <sys/pin.h>
```

```
kerrno_t pin_context_stack (flags)
```

```
long flags;
```

```
kerrno_t unpin_context_stack (flags)
```

```
long flags;
```

### Parameters

Item	Description
<i>flags</i>	Various flags to the kernel service. Must be set to 0.

### Description

Kernel code that pins its system call stack should call this service before the first kernel stack pin and call the `unpin_context_stack()` service after the last unpin. These services do not pin or unpin the C execution stack, but instead pin or unpin a hidden stack resource used for the kernel-key support.

### Execution Environment

These services must be called in the process environment.

### Return Values

Item	Description
0	Indicates a successful completion.
ENOMEM_PIN_CONTEXT_STACK	Indicates that the memory is not sufficient to satisfy the request.
ENOSPC_PIN_CONTEXT_STACK	Indicates that the page space is not sufficient.
EINVAL_PIN_CONTEXT_STACK	Indicates that the execution environment is not valid.
EINVAL_UNPIN_CONTEXT_STACK	Indicates that the execution environment is not valid. (For example, the service is not in the process environment or the kernel keys are not enabled or the value of the <i>flag</i> parameter is not valid.)

### Related reference:

“vm\_setseg\_kkey Kernel Service” on page 570

“vm\_protect\_kkey Kernel Service” on page 560

“raschk\_eaddr\_kkey Kernel Service” on page 436

“xmgethkeyset Kernel Service” on page 612

“xmsethkeyset Kernel Service” on page 613

## pinconf Kernel Service

### Purpose

Manages the list of free character buffers.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <cblock.h>
```

```
int pinconf ( delta)
int delta;
```

### Parameter

Item	Description
<i>delta</i>	Specifies the amount by which to change the number of free-pinned character buffers.

### Description

The **pinconf** service is used to control the size of the list of free-pinned character buffers. A positive value for the *delta* parameter increases the size of this list, while a negative value decreases the size.

All device drivers that use character blocks need to use the **pinconf** service. These drivers must indicate with a positive *delta* value the maximum number of character blocks they expect to be using concurrently. Device drivers typically call this service with a positive value when the **ddopen** routine is called. They should call the **pinconf** service with a negative value of the same amount when they no longer need the pinned character blocks. This occurs typically when the **ddclose** routine is called.

### Execution Environment

The **pinconf** kernel service can be called in the process environment only.

### Return Values

The **pinconf** service returns a value representing the amount by which the service changed the number of free-pinned character buffers.

#### Related reference:

“waitcfree Kernel Service” on page 583

#### Related information:

I/O Kernel Services

## pincode Kernel Service

### Purpose

Pins the code and data associated with a loaded object module.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/pin.h>
```

```
int pincode ( func)
int (*func) ();
```

## Parameter

Item	Description
<i>func</i>	Specifies an address used to determine the object module to be pinned. The address is typically that of a function exported by this object module.

## Description

The **pincode** service uses the **pin** service to pin the specified object module. The loader entry for the object module is used to determine the size of both the code and data.

## Execution Environment

The **pincode** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates successful completion.
EINVAL	Indicates that the <i>func</i> parameter is not a valid pointer to the function.
ENOMEM	Indicates that the <b>pincode</b> service was unable to pin the module due to insufficient real memory.

When an error occurs, the **pincode** service returns without pinning any pages.

### Related reference:

“pin Kernel Service” on page 408

### Related information:

Understanding Execution Environments

Memory Kernel Services

## pio\_assist Kernel Service

### Purpose

Provides a standardized programmed I/O exception handling mechanism for all routines performing programmed I/O.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
int pio_assist ( ioparms, iofunc, iorecov)
caddr_t ioparms;
int (*iofunc)( );
int (*iorecov)( );
```

### Parameters

Item	Description
<i>ioparms</i>	Points to parameters for the I/O routine.
<i>iofunc</i>	Specifies the I/O routine function pointer.
<i>iorecov</i>	Specifies the I/O recovery routine function pointer.

## Description

The **pio\_assist** kernel service assists in handling exceptions caused by programmed I/O. Use of the **pio\_assist** service standardizes the programmed I/O exception handling for all routines performing programmed I/O. The **pio\_assist** service is built upon other kernel services that routines access to provide their own exception handling if the **pio\_assist** service should not be used.

### Using the **pio\_assist** Kernel Service

To use the **pio\_assist** service, the device handler writer must provide a callable routine that performs the I/O operation. The device handler writer can also optionally provide a routine that can recover and log I/O errors. The mainline device handler code would then call the **pio\_assist** service with the following parameters:

- A pointer to the parameters needed by the I/O routine
- The function pointer for the routine performing I/O
- A pointer for the I/O recovery routine (or a null pointer, if there is no I/O recovery routine)

If the pointer for the I/O recovery routine is a null character, the *iofunc* routine is recalled to recover from I/O exceptions. The I/O routine for error retry should only be re-used if the I/O routine can handle being recalled when an error occurs, and if the sequence of I/O instructions can be reissued to recover from typical bus errors.

The *ioparms* parameter points to the parameters needed by the I/O routine. It is passed to the I/O routine when the **pio\_assist** service calls the I/O routine. It is also passed to the I/O recovery routine when the I/O recovery routine is invoked by the **pio\_assist** service. If any of the parameters found in the structure pointed to by the *ioparms* parameter are modified by the *iofunc* routine and needed by the *iorecov* or recalled *iofunc* routine, they must be declared as *volatile*.

### Requirements for Coding the Caller-Provided I/O Routine

The *iofunc* parameter is a function pointer to the routine performing the actual I/O. It is called by the **pio\_assist** service with the following parameters:

```
int iofunc (ioparms)
caddr_t ioparms;          /* pointer to parameters */
```

The *ioparms* parameter points to the parameters used by the I/O routine that was provided on the call to the **pio\_assist** kernel service.

If the **pio\_assist** kernel service is used with a null pointer to the *iorecov* I/O recovery routine, the *iofunc* I/O routine is called to retry all programmed I/O exceptions. This is useful for devices that have I/O operations that can be re-sent without concern for hardware state synchronization problems.

Upon return from the I/O, the return code should be 0 if no error was encountered by the I/O routine itself. If a nonzero return code is presented, it is used as the return code from the **pio\_assist** kernel service.

### Requirements for Coding the Caller-Provided I/O Recovery Routine

The *iorecov* parameter is a function pointer to the device handler's I/O recovery routine. This *iorecov* routine is responsible for logging error information, if required, and performing the necessary recovery operations to complete the I/O, if possible. This may in fact include calling the original I/O routine. The *iorecov* routine is called with the following parameters when an exception is detected during execution of the I/O routine:

```
int iorecov (parms, action, infop)
caddr_t parms; /* pointer to parameters passed to iofunc*/
int action;    /* action indicator */
struct pio_except *infop; /* pointer to exception info */
```

The *parms* parameter points to the parameters used by the I/O routine that were provided on the call to the **pio\_assist** service.

The *action* parameter is an operation code set by the **pio\_assist** kernel service to one of the following:

Item	Description
<b>PIO_RETRY</b>	Log error and retry I/O operations, if possible.
<b>PIO_NO_RETRY</b>	Log error but do not retry the I/O operation.

The **pio\_except** structure containing the exception information is platform-specific and defined in the **/usr/include/sys/except.h** file. The fields in this structure define the type of error that occurred, the bus address on which the error occurred, and additional platform-specific information to assist in the handling of the exception.

The *iorecov* routine should return with a return code of 0 if the exception is a type that the routine can handle. A **EXCEPT\_NOT\_HANDLED** return code signals that the exception is a type not handled by the *iorecov* routine. This return code causes the **pio\_assist** kernel service to invoke the next exception handler on the stack of exception handlers. Any other nonzero return code signals that the *iorecov* routine handled the exception but could not successfully recover the I/O. This error code is returned as the return code from the **pio\_assist** kernel service.

### Return Codes by the **pio\_assist** Kernel Service

The **pio\_assist** kernel service returns a return code of 0 if the *iofunc* I/O routine does not indicate any errors, or if programmed I/O exceptions did occur but were successfully handled by the *iorecov* I/O recovery routine. If an I/O exception occurs during execution of the *iofunc* or *iorecov* routines and the exception count has not exceeded the maximum value, the *iorecov* routine is called with an *op* value of **PIO\_RETRY**.

If the number of exceptions that occurred during this operation exceeds the maximum number of retries set by the platform-specific value of **PIO\_RETRY\_COUNT**, the **pio\_assist** kernel service calls the *iorecov* routine with an *op* value of **PIO\_NO\_RETRY**. This indicates that the I/O operation should not be retried. In this case, the **pio\_assist** service returns a return code value of **EIO** indicating failure of the I/O operation.

If the exception is not an I/O-related exception or if the *iorecov* routine returns with the return code of **EXCEPT\_NOT\_HANDLED** (indicating that it could not handle the exception), the **pio\_assist** kernel service does not return to the caller. Instead, it invokes the next exception handler on the stack of exception handlers for the current process or interrupt handler. If no other exception handlers are on the stack, the default exception handler is invoked. The normal action of the default exception handler is to cause a system crash.

### Execution Environment

The **pio\_assist** kernel service can be called from either the process or interrupt environment.

## Return Values

Item	Description
0	Indicates that either no errors were encountered, or PIO errors were encountered and successfully handled.
EIO	Indicates that the I/O operation was unsuccessful because the maximum number of I/O retry operations was exceeded.

### Related information:

Kernel Extension and Device Driver Management Kernel Services

User-Mode Exception Handling

## Process State-Change Notification Routine

### Purpose

Allows kernel extensions to be notified of major process and thread state transitions.

### Syntax

```
void prochadd_handler ( term, type, id)
struct proch *term;
int type;
long id;
```

```
void proch_reg_handler ( term, type, id)
struct prochr *term;
int type;
long id;
```

### Parameters

Item	Description
<i>term</i>	Points to the <b>proch</b> structure used in the <b>prochadd</b> call or to the <b>prochr</b> structure used in the <b>proch_reg</b> call.

<b>Item</b>	<b>Description</b>
<i>type</i>	<p>Defines the state change event being reported: process initialization, process termination, process exec, thread initialization, or thread termination. These values are defined in the <code>/usr/include/sys/proc.h</code> file. The values that may be passed as <i>type</i> also depend on how the callout is requested.</p> <p>Possible <code>prochadd_handler</code> <i>type</i> values:</p> <p><b>PROCH_INITIALIZE</b> Process is initializing.</p> <p><b>PROCH_TERMINATE</b> Process is terminating.</p> <p><b>PROCH_EXEC</b> Process is about to exec a new program.</p> <p><b>THREAD_INITIALIZE</b> A new thread is created.</p> <p><b>THREAD_TERMINATE</b> A thread is terminated.</p> <p>Possible <code>prochr_reg_handler</code> <i>type</i> values:</p> <p><b>PROCHR_INITIALIZE</b> Process is initializing.</p> <p><b>PROCHR_TERMINATE</b> Process is terminating.</p> <p><b>PROCHR_EXEC</b> Process is about to exec a new program.</p> <p><b>PROCHR_THREAD_INIT</b> A new thread is created.</p> <p><b>PROCHR_THREAD_TERM</b> A thread is terminated.</p>
<i>id</i>	Defines either the process ID or the thread ID.

## Description

The notification callout is set up by using either the **prochadd** or the **prochr\_reg** kernel service. If you request the notification using the **prochadd** kernel service, the callout follows the syntax shown first as **prochadd\_handler**. If you request the notification using the **prochr\_reg** kernel service, the callout follows the syntax shown second as **prochr\_reg\_handler**.

For process initialization, the **process state-change notification** routine is called in the execution environment of a parent process for the initialization of a newly created child process. For kernel processes, the notification routine is called when the **initp** kernel service is called to complete initialization.

For process termination, the notification routines are called before the kernel handles default termination procedures. The routines must be written so as not to allocate any resources under the terminating process. The notification routine is called under the process image of the terminating process.

### Related reference:

“prochadd Kernel Service” on page 417

“prochdel Kernel Service” on page 419

### Related information:

Kernel Extension and Device Driver Management Kernel Services

## prochr\_reg Kernel Service

### Purpose

Registers a callout handler.

### Syntax

```
#include <sys/proc.h>
```

```
int prochr_reg(struct prochr *)
```

**Note:** The `prochr` structure contains the following elements that must be set prior to calling `prochr_reg`:

```
void (* prochr_handler)(struct prochr *, int, long)
unsigned int int prochr_mask
```

### Parameters

#### Item

*int prochr\_mask*

#### Description

Specifies the set of kernel events for which a callout is requested. Unlike the `old_style` interface, the callout is invoked only for the specified events. This mask is formed by ORing together any of these defined values:

**PROCHR\_INITIALIZE**  
Process created.

**PROCHR\_TERMINATE**  
Process terminated

**PROCHR\_EXEC**  
Process has issued the `exec` system call

**PROCHR\_THREADINIT**  
Thread created

**PROCHR\_THREADTERM**  
Thread terminated

*prochr\_handler*

Specifies the callout function to be called when specified kernel events occur.

### Description

If the same `struct prochr *` is registered more than once, only the most recently specified information is retained in the kernel.

The `struct prochr *` is not copied to a new location in memory. As a result, if the structure is changed, results are unpredictable. This structure does not need to be pinned.

The primary consideration for the new-style interface is to improve scalability. A lock is only acquired when callouts are made. A summary mask of all currently registered callout event types is maintained. This summary mask is updated every time `prochr_reg` or `prochr_unreg` is called, even when registering an identical `struct prochr *`. Further, the lock is a complex lock, so once callouts have been registered, there is no lock contention in invoking them because the lock is held read-only.

When a callout to a registered handler function is made, the parameters passed are:

- a pointer to the registered `prochr` structure
- a callout request value to indicate the reason for the callout
- a thread or process ID



## Return Values

On successful completion, the **proch\_reg** kernel service returns a value of 0. The only error (non-zero) return is from trying to register with a NULL pointer.

## Execution Environment

The **proch\_reg** kernel service can be called from the process environment only.

### Related reference:

“proch\_unreg Kernel Service”

“Process State-Change Notification Routine” on page 414

### Related information:

Kernel Extension and Driver Management Kernel Services

## proch\_unreg Kernel Service

### Purpose

Unregisters a callout handler that was previously registered using the **proch\_reg** kernel service.

### Syntax

```
#include <sys/proc.h>
int proch_unreg(struct prochr *old_prochr);
```

### Parameter

Item	Description
<i>old_prochr</i>	Specifies the address of the <b>proch</b> structure to be unregistered.

### Description

Unregisters an existing callout handler that was previously registered using the **proch\_reg(0)** kernel service.

## Return Values

On successful completion, the **proch\_unreg** kernel service returns a value of 0. An error (non-zero) return occurs when trying to unregister a handler that is not presently registered.

## Execution Environment

The **proch\_unreg** kernel service can be called from the process environment only.

### Related reference:

“proch\_reg Kernel Service” on page 416

### Related information:

Kernel Extension and Driver Management Kernel Services

## prochadd Kernel Service

### Purpose

Adds a system-wide process state-change notification routine.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/proc.h>
```

```
void prochadd ( term)
struct proch *term;
```

## Parameters

Item	Description
<i>term</i>	Points to a <b>proch</b> structure containing a notification routine to be added from the chain of systemwide notification routines.

## Description

The **prochadd** kernel service allows kernel extensions to register for notification of major process state transitions. The **prochadd** service allows the caller to be notified when a process:

- Has just been created.
- Is about to be terminated.
- Is executing a new program.

The complete list of callouts is:

Callout	Description
PROCH_INITIALIZE	Process (pid) created ( <b>initp</b> , <b>kforkx</b> )
PROCH_TERMINATE	Process (pid) terminated ( <b>kexitx</b> )
PROCH_EXEC	Process (pid) executing ( <b>execvex</b> )
THREAD_INITIALIZE	Thread (tid) created ( <b>kforkx</b> , <b>thread_create</b> )
THREAD_TERMINATE	Thread (tid) created ( <b>kexitx</b> , <b>thread_terminate</b> )

The **prochadd** service is typically used to allow recovery or reassignment of resources when processes undergo major state changes.

The caller should allocate a **proch** structure and update the **proch.handler** field with the entry point of a caller-supplied notification routine before calling the **prochadd** kernel service. This notification routine is called once for each process in the system undergoing a major state change.

The **proch** structure has the following form:

```
struct proch
{
    struct proch *next
    void          *handler ();
}
```

## Execution Environment

The **prochadd** kernel service can be called from the process environment only.

**Related reference:**

“prochdel Kernel Service” on page 419

“Process State-Change Notification Routine” on page 414

**Related information:**

Kernel Extension and Driver Management Kernel Services

## prochdel Kernel Service

### Purpose

Deletes a process state change notification routine.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/proc.h>
```

```
void prochdel ( term)
struct proch *term;
```

### Parameter

Item	Description
<i>term</i>	Points to a <b>proch</b> structure containing a notification routine to be removed from the chain of system-wide notification routines. This structure was previously registered by using the <b>prochadd</b> kernel service.

### Description

The **prochdel** kernel service removes a process change notification routine from the chain of system-wide notification routines. The registered notification routine defined by the handler field in the **proch** structure is no longer to be called by the kernel when major process state changes occur.

If the **proch** structure pointed to by the *term* parameter is not found in the chain of structures, the **prochdel** service performs no operation.

### Execution Environment

The **prochdel** kernel service can be called from the process environment only.

#### Related reference:

“prochadd Kernel Service” on page 417

“Process State-Change Notification Routine” on page 414

#### Related information:

Kernel Extension and Driver Management Kernel Services

## probe or kprobe Kernel Service

### Purpose

Logs errors with symptom strings.

### Library (for probe)

Run-time Services Library.

### Syntax

```
#include <sys/probe.h>
or
#include <sys/sysprobe.h>
int probe ( probe_p)
probe_t *probe_p
int kprobe (probe_p)
probe_t *probe_p
```

## Description

The `probe` subroutine logs an entry to the error log. The entry consists of an error log entry as defined in the `errlog` subroutine and the `err_rec.h` header file, and a symptom string.

The `probe` subroutine is called from an application, while `kprobe` is called from the Kernel and Kernel extensions. Both `probe` and `kprobe` have the same interfaces, except for return codes.

IBM software should use the `sys/sysprobe.h` header file while non-IBM programs should include the `sys/probe.h` file. This is because IBM symptom strings must conform to different rules than non-IBM strings. It also tells any electronic support application whether or not to route the symptom string to IBM's Retain database.

## Parameters

Item	Description
<code>probe_p</code>	<p>is a pointer to the data structure which contains the pointer and length of the error record, and the data for the probe. The error record is described under the <code>errlog</code> subroutine and defined in <code>err_rec.h</code>.</p> <p>The first word of the structure is a magic number to identify this version of the structure. The magic number should be set to <code>PROBE_MAGIC</code>.</p> <p><b>Note:</b> <code>PROBE_MAGIC</code> is different between <code>probe.h</code> and <code>sysprobe.h</code> to distinguish an IBM symptom string from a non-IBM string.</p> <p>The probe data consists of flags which control probe handling, the number of symptom string keywords, followed by an array consisting of one element for each keyword.</p>

## Flags

Item	Description
<code>SSNOSEND</code>	indicates this symptom string shouldn't be forwarded to automatic problem opening facilities. An example where <code>SSNOSEND</code> should be used is in symptom data used for debugging purposes.
<code>nsskwd</code>	This gives the number of keywords specified (i.e.), the number of elements in the <code>sskwd</code> s array.

**Item**  
**sskwd**

**Description**

This is an array of keyword/value pairs. The keywords and their values are in the following table. The *I/S* value indicates whether the *keyword* and *value* are informational or are part of the logged symptom string. The number in parenthesis indicates, where applicable, the maximum string length.

<b>keyword</b>	<b>I/S</b>	<b>value</b>	<b>type</b>	<b>Description</b>
SSKWD_LONGNAME	I	char *	(30)	Product's long name
SSKWD_OWNER	I	char *	(16)	Product's owner
SSKWD_PIDS	S	char *	(11)	product id.
(required for IBM symptom strings)				
SSKWD_LVL	S	char *	(5)	product level
(required for IBM symptom strings)				
SSKWD_APPLID	I	char *	(8)	application id.
SSKWD_PCCS	S	char *	(8)	probe id
(required for all symptom strings)				
SSKWD_DESC	I	char *	(80)	problem description
SSKWD_SEV	I	int		severity from
				1 (highest) to 4 (lowest).
				3 is the default.
SSKWD_AB	S	char *	(5)	abend code
SSKWD_ADRS	S	void *		address. If used at all,
				this should be a relative address.
SSKWD_DEVS	S	char *	(6)	Device type
SSKWD_FLDS	S	char *	(9)	arbitrary character string.
				This is usually a field name and
				the SSKWD_VALUE
				keyword specifies the value.
SSKWD_MS	S	char *	(11)	Message number
SSKWD_OPCS	S	char *	(8)	OP code
SSKWD_OVS	S	char *	(9)	overwritten storage
SSKWD_PRC	S			unsigned long return code
SSKWD_REGS	S	char *	(4)	Register name (e.g.)
				GR15 or LR unsigned long Value
SSKWD_VALU	S			
SSKWD_RIDS	S	char *	(8)	resource or module id.
SSKWD_SIG	S	int		Signal number
SSKWD_SN	S	char *	(7)	Serial Number
SSKWD_SRN	S	char *	(9)	Service Req. Number If specified,
				and no error is logged,
				a hardware error is assumed.
SSKWD_WS	S	char *	(10)	Coded wait

**Note:** The **SSKWD\_PCCS** value is always required. This is the probe id. Additionally, for IBM symptom strings, the **SSKWD\_PIDS** and **SSKWD\_LVL** keywords are also required

If either the **erecp** or **erecl** fields in the **probe\_rec** structure is 0 then no error logging record is being passed, and one of the default templates for symptom strings is used. The default template indicating a software error is used unless the **SSKWD\_SRN** keyword is specified. If it is, the error is assumed to be a hardware error. If you don't want to log your own error with a symptom string, and you want to have a hardware error, and don't want to use the **SSKWD\_SRN** value, then you can supply an error log record using the error identifier of **ERRID\_HARDWARE\_SYMPTOM**, see the **/usr/include/sys/errids.h** file.

**Return Values for probe Subroutine**

Item	Description
0	Successful
-1	Error. The errno variable is set to
EINVAL	Indicates an invalid parameter
EFAULT	Indicates an invalid address

## Return Values for kprobe Kernal Service

Item	Description
0	Successful
EINVAL	Indicates an invalid parameter

## Execution Environment

**probe** is executed from the application environment.

**kprobe** is executed from the Kernel and Kernel extensions. Currently, **kprobe** must not be called with interrupts disabled.

## Files

Item	Description
<code>/usr/include/sys/probe.h</code>	Contains parameter definition.

### Related reference:

“errsave or errlast Kernel Service” on page 139

### Related information:

Error Logging Overview  
errlog subroutine

## purblk Kernel Service

### Purpose

Purges the specified block from the buffer cache.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/buf.h>
```

```
void purblk ( dev, blkno)
dev_t dev;
daddr_t blkno;
```

### Parameters

Item	Description
<i>dev</i>	Specifies the device containing the block to be purged.
<i>blkno</i>	Specifies the block to be purged.

## Description

The **purblk** kernel service purges (that is, makes unreclaimable by marking the block with a value of **STALE**) the specified block from the buffer cache.

## Execution Environment

The **purblk** kernel service can be called from the process environment only.

## Return Values

The **purblk** service has no return values.

### Related reference:

“brelse Kernel Service” on page 32

“geteblk Kernel Service” on page 186

### Related information:

Block I/O Buffer Cache Kernel Services: Overview

## putc Kernel Service

### Purpose

Places a character at the end of a character list.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <cblock.h>
```

```
int putc ( c, header)
char c;
struct clist *header;
```

### Parameters

Item	Description
<i>c</i>	Specifies the character to place on the character list.
<i>header</i>	Specifies the address of the <b>clist</b> structure that describes the character list.

## Description

**Attention:** The caller of the **putc** service must ensure that the character list is pinned. This includes the **clist** header and all the **cblock** character buffers. Character blocks acquired from the **getc** service are also pinned. Otherwise, the system may crash.

The **putc** kernel service puts the character specified by the *c* parameter at the end of the character list pointed to by the *header* parameter.

If the **putc** service indicates that there are no more buffers available, the **waitcfree** service can be used to wait until a character block is available.

## Execution Environment

The **putc** kernel service can be called from either the process or interrupt environment.

## Return Values

Item	Description
0	Indicates successful completion.
-1	Indicates that the character list is full and no more buffers are available.

### Related reference:

“**getc** Kernel Service” on page 183

“**putc** Kernel Service” on page 426

### Related information:

I/O Kernel Services

## putcb Kernel Service

### Purpose

Places a character buffer at the end of a character list.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <cblock.h>
```

```
void putcb ( p, header)
struct cblock *p;
struct clist *header;
```

### Parameters

Item	Description
<i>p</i>	Specifies the address of the character buffer to place on the character list.
<i>header</i>	Specifies the address of the <b>clist</b> structure that describes the character list.

### Description

**Attention:** The caller of the **putcb** service must ensure that the character list is pinned. This includes the **clist** header and all the **cblock** character buffers. Character blocks acquired from the **getc** service are pinned. Otherwise, the system may crash.

The **putcb** kernel service places the character buffer pointed to by the *p* parameter on the end of the character list specified by the *header* parameter. Before calling the **putcb** service, you must load this new buffer with characters and set the *c\_first* and *c\_last* fields in the **cblock** structure. The *p* parameter is the address returned by either the **getc** or the **getc** service.

## Execution Environment

The **putcb** kernel service can be called from either the process or interrupt environment.

## Return Values



Item	Description
0	Indicates successful completion.
-1	Indicates that the character list is full and no more buffers are available.

#### Related reference:

“getcb Kernel Service” on page 183

“putcf Kernel Service” on page 426

#### Related information:

I/O Kernel Services

## putcbp Kernel Service

### Purpose

Places several characters at the end of a character list.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <cblock.h>
```

```
int putcbp ( header, source, n)
struct clist *header;
char *source;
int n;
```

### Parameters

Item	Description
<i>header</i>	Specifies the address of the <b>clist</b> structure that describes the character list.
<i>source</i>	Specifies the address from which characters are read to be placed on the character list.
<i>n</i>	Specifies the number of characters to be placed on the character list.

### Description

**Attention:** The caller of the **putcbp** service must ensure that the character list is pinned. This includes the **clist** header and all of the **cblock** character buffers. Character blocks acquired from the **getc** service are pinned. Otherwise, the system may crash.

The **putcbp** kernel service operates on the characters specified by the *n* parameter starting at the address pointed to by the *source* parameter. This service places these characters at the end of the character list pointed to by the *header* parameter. The **putcbp** service then returns the number of characters added to the character list. If the character list is full and no more buffers are available, the **putcbp** service returns a 0. Otherwise, it returns the number of characters written.

### Execution Environment

The **putcbp** kernel service can be called from either the process or interrupt environment.

### Return Values

The **putcbp** service returns the number of characters written or a value of 0 if the character list is full, and no more buffers are available.

#### Related reference:

“pincl Kernel Service” on page 410

“putcf Kernel Service”

“waitcfree Kernel Service” on page 583

## **putcf Kernel Service**

### **Purpose**

Frees a specified buffer.

### **Syntax**

```
#include <sys/types.h>
#include <sys/errno.h>
#include <cblock.h>
```

```
void putcf ( p )
struct cblock *p;
```

### **Parameter**

Item	Description
<i>p</i>	Identifies which character buffer to free.

### **Description**

The **putcf** kernel service unpins the indicated character buffer.

The **putcf** service returns the specified buffer to the list of free character buffers.

### **Execution Environment**

The **putcf** kernel service can be called from either the process or interrupt environment.

### **Return Values**

The **putcf** service has no return values.

#### **Related information:**

I/O Kernel Services

## **putcfl Kernel Service**

### **Purpose**

Frees the specified list of buffers.

### **Syntax**

```
#include <sys/types.h>
#include <sys/errno.h>
#include <cblock.h>
```

```
void putcfl ( header )
struct clist *header;
```

### **Parameter**

Item	Description
<i>header</i>	Identifies which list of character buffers to free.

## Description

The **putcfl** kernel service returns the specified list of buffers to the list of free character buffers. The **putcfl** service unpins the indicated character buffer.

**Note:** The caller of the **putcfl** service must ensure that the header and **clist** structure are pinned.

## Execution Environment

The **putcfl** kernel service can be called from either the process or interrupt environment.

## Return Values

The **putcfl** service has no return values.

### Related information:

I/O Kernel Services

## putcx Kernel Service Purpose

Places a character on a character list.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/cblock.h>
```

```
int putcx ( c, header)
char c;
struct clist *header;
```

## Parameters

Item	Description
<i>c</i>	Specifies the character to place at the front of the character list.
<i>header</i>	Specifies the address of the <b>clist</b> structure that describes the character list.

## Description

The **putcx** kernel service puts the character specified by the *c* parameter at the front of the character list pointed to by the *header* parameter. The **putcx** service is identical to the **putc** service, except that it puts the character at the front of the list instead of at the end.

If the **putcx** service indicates that there are no more buffers available, the **waitcfree** service can be used to wait until a character buffer is available.

**Note:** The caller of the **putcx** service must ensure that the character list is pinned. This includes the **clist** header and all the **cblock** character buffers. Character blocks acquired from the **getcfl** service are pinned.

## Execution Environment

The **putcx** kernel service can be called from either the process or interrupt environment.

## Return Values

Item	Description
0	Indicates successful completion.
-1	Indicates that the character list is full and no more buffers are available.

### Related reference:

“pincf Kernel Service” on page 410

“putcfl Kernel Service” on page 426

### Related information:

I/O Kernel Services

## q

The following kernel services begin with the with the letter q.

### query\_proc\_info Kernel Service

#### Purpose

Returns specific information about the current process or thread.

#### Syntax

```
#include <sys/encap.h>
```

```
int query_proc_info (type)
int type;
```

#### Parameters

Item	Description
<i>type</i>	Specifies the type of process or thread information requested. The <i>type</i> parameter can be one of the following values:  <b>QPI_XPG_SUS_ENV</b> Queries whether the calling process has SPEC 1170 environment active.  <b>QTI_FUNNELLED</b> Queries whether the current thread is funneled.

#### Description

The **query\_proc\_info** kernel service returns information about the current process or thread.

When called with the value **QPI\_XPG\_SUS\_ENV** as the *type* parameter, it returns TRUE (1) when the process has SPEC 1170 active, that is, the process was issued with the environment variable **XPG\_SUS\_ENV** defined. Otherwise, the routine returns FALSE (0). When called with the value **QTI\_FUNNELLED** as the *type* parameter, the **query\_proc\_info** kernel service returns TRUE (1) if the current thread has been funneled. Otherwise, the routine returns FALSE (0).

#### Execution Environment

The **query\_proc\_info** kernel service can be called from either the process or interrupt environment.

#### Return Values

Item	Description
1	True.
0	False.

## r

The following kernel services begin with the with the letter r.

### **RAS\_BLOCK\_NULL Exported Data Structure Purpose**

Allows for the silent failure of **ras\_register** calls due to memory allocation errors.

#### **Syntax**

```
#include <sys/ras.h>
```

```
extern const ras_block_t RAS_BLOCK_NULL
```

#### **Description**

The **RAS\_BLOCK\_NULL** data structure allows components to go through their normal code paths when they receive an ENOMEM error from the **ras\_register** kernel service. The presence of this data structure does not need to be explicitly checked by callers of RAS functions. All RAS domain functions (such as Component Tracing) are disabled with this control block.

#### **Related reference:**

“**ras\_register** and **ras\_unregister** Exported Kernel Services” on page 432

“**ras\_customize** Exported Kernel Service” on page 430

#### **Related information:**

CT\_HOOKx subroutine

### **ras\_control Exported Kernel Service Purpose**

Controls component RAS characteristics.

#### **Syntax**

```
#include <sys/ras.h>
```

```
kernno_t ras_control (
ras_block_t ras_blk,
ras_cmd_t command,
void * arg,
long argsize);
```

#### **Description**

The **ras\_control** kernel service passes a command to the callback for the component referenced by the **ras\_blk** parameter. If the **ras\_blk** parameter is not known, use the **ras\_path\_control** call.

**Note:** During the **ras\_control** process, callbacks to the registrant of the component might be initiated for changes that the RAS infrastructure makes to the component. The registrant should be aware of this for locking purposes (for instance, the registrant should not hold any locks that the callback needs).

If the *ras\_blk* input parameter has a value of `RAS_BLOCK_NULL`, the **ras\_control** kernel service returns without errors and takes no action.

## Parameters

Item	Description
<i>ras_blk</i>	The target control block pointer.
<i>command</i>	Command passed to the callback. Commands are specific to a given RAS domain, such as Component Trace.
<i>arg</i>	Optional argument for the command.
<i>argsize</i>	Size of the argument, if a buffer or structure.

## Execution Environment

The calling environment of the **ras\_control** kernel service varies by individual command. The calling environment of a particular command is documented with the command itself.

## Return Values

The **ras\_control** kernel service returns 0 for success and a non-zero error code for failure.

### Related reference:

“*ras\_customize* Exported Kernel Service”

“*ras\_path\_control* Exported Kernel Services” on page 431

### Related information:

Component Trace Facility

*ras\_callback* subroutine

## **ras\_customize** Exported Kernel Service Purpose

Loads persistent customized properties for a RAS control block.

## Syntax

```
#include <sys/ras.h>
```

```
kernno_t ras_customize (ras_block_t ras_blk);
```

## Description

The **ras\_customize** kernel service checks for, and applies persistent customized properties for a given *ras\_blk* parameter. After applying any persistent properties, the **ras\_customize** kernel service puts the *ras\_blk* parameter in a usable state. Registration is not complete without a call to the **ras\_customize** kernel service.

**Note:** During the **ras\_customize** process, callbacks to the registrant might be initiated for changes that the RAS infrastructure makes to the component. The registrant should be aware of this for locking and initialization purposes (for example, the registrant should not be holding any locks that the callback needs, and the private data for the callback should be initialized before **ras\_customize** is called).

If the *ras\_blk* input parameter has a value of `RAS_BLOCK_NULL`, the **ras\_customize** kernel service returns without errors and takes no action.

## Parameters

Item	Description
<i>ras_blk</i>	The control block to act on. Must be previously allocated by the <b>ras_register</b> kernel service.

## Execution Environment

The **ras\_customize** kernel service must be called from the process environment.

## Return Values

Item	Description
0	Successful.
non-zero	Unsuccessful.

### Related reference:

“ras\_control Exported Kernel Service” on page 429

### Related information:

Component Trace Facility

*ras\_callback* subroutine

## ras\_path\_control Exported Kernel Services Purpose

Controls component RAS characteristics.

## Syntax

```
#include <sys/ras.h>
```

```
kernno_t ras_path_control (
char * path,
ras_cmd_t command,
void * arg,
long argsize);
```

## Description

The **ras\_path\_control** kernel service passes a command to the RAS component specified by the *path* parameter.

**Note:** During the **ras\_path\_control** process, callbacks to the registrant of the component might be initiated for changes that the RAS infrastructure makes to the component. The registrant should be aware of this for locking purposes (for instance, the registrant should not be holding any locks the callback needs).

## Parameters

Item	Description
<i>path</i>	The pathname of the component to receive the <i>command</i> parameter.
<i>command</i>	Command passed to the callback. Commands are specific to a given RAS domain, such as Component Trace.
<i>arg</i>	Optional argument for the command.
<i>argsize</i>	Size of the argument, if a buffer or structure.

## Execution Environment

The calling environment of the **ras\_path\_control** kernel service varies by individual command. The calling environment of a particular command is documented with the command itself.

## Return Values

Item	Description
0	Successful.
non-zero	Unsuccessful.

### Related reference:

“**ras\_control** Exported Kernel Service” on page 429

“**ras\_customize** Exported Kernel Service” on page 430

### Related information:

Component Trace Facility

## **ras\_register** and **ras\_unregister** Exported Kernel Services Purpose

Registers and unregisters a RAS component.

### Syntax

```
#include <sys/ras.h>
```

```
kernno_t ras_register (  
ras_block_t * rasbp,  
char * name,  
ras_block_t parent,  
ras_type_t typesubtype,  
char * desc,  
long flags,  
ras_callback_t ras_callback,  
void * private_data);
```

```
kernno_t ras_unregister (ras_block_t ras_blk);
```

### Description

The **ras\_register** kernel service and the **ras\_unregister** kernel service register and unregister RAS handlers which are invoked by the kernel when the system needs to communicate various RAS commands to each component.

The **ras\_register** kernel service registers a component with the given name under the *parent* provided. If the parent is NULL, the **ras\_register** kernel service registers name as a base component, but the *typesubtype* parameter must be provided. The *name* parameter specifies the name for the subcomponent or base component (it is not a full component path). The *flags* field is used to specify what aspects of RAS the component understands. The *ras\_callback* is the mechanism by which the RAS subsystem communicates various commands to the component, depending on what aspects of RAS the component understands. The *desc* parameter provides a short description for the component as a service aid.

The **ras\_register** kernel service allocates a *ras\_block\_t* member and returns the control block for the component through the *rasbp* argument. This control block can be used in **ras\_control** calls and further **ras\_register** calls (to allocate children, for instance).



If the registration fails due to the system being out of memory, the value of the *rasbp* argument is set to `RAS_BLOCK_NULL`. All RAS functions for this component are disabled. RAS kernel services accept `RAS_BLOCK_NULL` control blocks but take no action. If the control block is set to `RAS_BLOCK_NULLRAS`, domain related functions (such as the `CT_HOOKx` and `CT_GEN` macros) run correctly but take no action. This action allows the `ENOMEM` type failures from the `ras_register` kernel service to be safely ignored. The value of the *rasbp* argument for all other types of errors is undefined.

The `ras_unregister` kernel service unregisters a component previously registered with the `ras_register` kernel service. The *ras\_blk* parameter should have no further children.

## Parameters

Item	Description
<i>rasbp</i>	The newly allocated <code>ras_block_t</code> member.
<i>name</i>	The name of the component, not its full pathname. Individual node names are limited to the number of characters specified by the value of the <code>RAS_NAME_MAX</code> parameter (including the terminating NULL character). The full component path (the concatenated names of a child component and all of its ancestors) is limited to the number of characters specified by the value of the <code>RAS_PATH_MAX</code> parameter (including the terminating NULL character). The <code>ras_register</code> kernel service reconstructs the full component path and rejects registrations for components whose full path exceeds the value of the <code>RAS_PATH_MAX</code> parameter. Node names are restricted to the character set "A-Z", "a-z", "0-9" and "_".
<i>parent</i>	An optional pointer to the parent component or NULL if none.
<i>typesubtype</i>	If parent is NULL, mandatory parameter is used to categorize the component. The top 16-bits of the lower word of this field are the type, and the bottom 16-bits are the subtype. The <i>typesubtype</i> is a <code>ras_type_t</code> member, which is an enum. See the <code>sys/ras_base.h</code> file for a description of the types available. If parent is non-NULL, this parameter is required to be the value of the <code>RAS_TYPE_CHILD</code> parameter.
<i>desc</i>	A short description string for the component. The <i>desc</i> string is limited to the number of characters specified by the value of the <code>RAS_DESC_MAX</code> parameter (including the terminating null). The <i>desc</i> string has no character set restriction. Any static elements of the string should be in U.S. English, but dynamic elements have no restriction.
<i>flags</i>	Indicates what type of RAS systems this component is aware of. Valid choices are the following: <ul style="list-style-type: none"> <li>• <code>RASF_TRACE_AWARE</code>: Component is Component Trace aware.</li> <li>• <code>RASF_ERROR_AWARE</code>: Component is Error Checking aware.</li> </ul> These flags are defined in the <code>sys/ras.h</code> file.
<i>ras_callback</i>	A function pointer provided by the registrant and called by the framework each time an external event modifies a property of the component. See the <code>ras_callback</code> interface specification.
<i>private_data</i>	An optional pointer to a component private memory area passed to the <code>ras_callback</code> function upon callback.
<i>ras_blk</i>	The control block to remove.

## Execution Environment

Both the `ras_register` kernel service and the `ras_unregister` kernel service must be called from the process environment.

## Return Values

The following are the return values of the `ras_register` kernel service.

Item	Description
0	Successful.
non-zero	Unsuccessful.

The following are the return values of the **ras\_unregister** kernel service.

Item	Description
0	Successful.
non-zero	Unsuccessful.

#### Related reference:

“ras\_customize Exported Kernel Service” on page 430

#### Related information:

Component Trace Facility

ras\_callback subroutine

## ras\_ret\_query\_parms Kernel Service

### Purpose

Returns callback parameters in the **ras\_query\_parms** structure.

### Syntax

```
#include <sys/ras.h>
```

```
kernno_t ras_ret_query_parms (retp, fmtstr, numstrings, descr)
```

```
ras_query_parms_t *retp;
```

```
char *fmtstr;
```

```
int numstrings;
```

```
char *descr[];
```

### Parameters

Item	Description
<i>retp</i>	Points to the <b>ras_query_parms_t</b> data item to be filled in.
<i>fmtstr</i>	This is a format specifier. It has the following form: <i>spec-list</i>  or kywd=spec-list kywd=spec-list ...  Where the <i>spec-list</i> variable is of the form: <i>spec,spec,...</i> . The <i>spec</i> variable must be %x, %xx, %d, %dd, %s, or %ss. If the characters x, d, or s are doubled, for example, %xx, this indicates that multiple values are allowed.  The following are some valid <i>fmtstr</i> values:  %x        One hexadecimal value.  %x,%d    One hexadecimal and one decimal value.  %xx      Multiple hexadecimal values.  k1=%x,%d k2=%dd Keyword k1 takes one hexadecimal value and one decimal value. Keyword k2 takes multiple decimal values.
<i>numstrings</i>	Specifies the number of strings in the <i>descr</i> string array. The value must be at least 1.

Item	Description
<i>descr</i>	<p>Specifies the component and parameters. There must be at least one string. The first string describes the component's function. If the component takes positional parameters, the following string(s) describe those. If keyword parameters are supplied, each keyword must have a corresponding <i>descr</i> string in the array describing that keyword.</p> <p>The <code>ras_ret_query_parms</code> kernel service does not return an error if the number of the <i>descr</i> strings does not match the format string. Instead, either the last keywords do not have help text, or the excess help strings are simply displayed.</p>

## Description

The `ras_ret_query_parms` kernel service can be used by a callback to aid in filling in the `ras_query_parms_t` structure when it receives the `RASC_QUERY_PARMS` call. This function formats the help text and places it into the `ras_query_parms_t` structure. If there is insufficient space for the help text in the provided `ras_query_parms_t` item, it returns `ENOMEM_RASC_CONTROL_QUERYPARMS`. The callback then just returns this error code.

The help text provided must follow the following conventions:

*component* - first line of description  
*component:parameters* - parameter(s) description

or

*component* - first line of description  
*component:kywd1=parms* - *kywd1:parms* description  
*component:kywd2=parms* - *kywd2:parms* description

## Execution Environment

The `ras_ret_query_parms` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful completion.
<code>EINVAL_RASC_CONTROL_QUERYPARMS</code>	Indicates that one or more parameters was not valid.
<code>EFAULT_RASC_CONTROL_QUERYPARMS</code>	Indicates that one or more parameter addresses was not valid.
<code>ENOMEM_RASC_CONTROL_QUERYPARMS</code>	Indicates that the <code>rqp_text</code> size was not large enough.

### Related reference:

“`dmp_compspec` and `dmp_compext` Kernel Services” on page 93

## raschk\_eaddr\_hkeyset Kernel Service

### Purpose

Checks if an effective address can be referenced with a hardware keyset.

### Syntax

```
#include <sys/raschk.h>
#include <sys/skeys.h>
#include <sys/kerrno.h>
```

```
kerrno_t raschk_eaddr_hkeyset (eaddr, hset, flags)
void * eaddr;
hkeyset_t hset;
unsigned long flags;
```

## Parameters

Item	Description
<i>eaddr</i>	Effective address to validate. Only one byte is checked.
<i>hset</i>	Hardware keyset to validate against.
<i>flags</i>	The following flags are defined: <b>RCHK_EHK_NOFAULT</b> No page faults are permitted while performing this check. <b>RCHK_EHK_NOPAGEIN</b> No page in is performed during this check. <b>RCHK_EHK_READ</b> Validates for read access. <b>RCHK_EHK_WRITE</b> Validates for write access.

## Description

The `raschk_eaddr_hkeyset` kernel service performs an advisory runtime check to determine if an effective address can be referenced with a hardware keyset.

Read and write access checks are independently specified in the *flags* field. A check for read and write access requires both flags to be set.

## Execution Environment

The `raschk_eaddr_hkeyset` kernel service can be called from the process or interrupt environment.

## Return Values

Item	Description
0	Successful.
EFAULT_RASCHK_EADDR_HKEYSET	Operation failed because a page in or page fault was not allowed.
EFAULT_RASCHK_EADDR_HKEYSET_PROT	The address failed the protection check.
EINVAL_RASCHK_EADDR_HKEYSET	The address to validate was determined to be invalid, or neither READ nor WRITE checking was requested.

### Related reference:

“`raschk_eaddr_kkey` Kernel Service”

## `raschk_eaddr_kkey` Kernel Service

### Purpose

Checks if an effective address can be referenced with a kernel-key.

### Syntax

```
#include <sys/raschk.h>  
#include <sys/kerrno.h>
```

```
kerrno_t raschk_eaddr_kkey (eaddr, kkey, flags)  
void * eaddr;  
kkey_t kkey;  
unsigned long flags;
```

## Parameters

Item	Description
<i>eaddr</i>	Effective address to validate. Only one byte is checked.
<i>kkey</i>	Kernel-key to check.
<i>flags</i>	The following flags are defined:
	<b>RCHK_EK_NOFAULT</b> No page faults of any kind are permitted while performing this check.
	<b>RCHK_EK_NOPAGEIN</b> No page in will be performed during this check.

## Description

The `raschk_eaddr_kkey` kernel service performs an advisory runtime check to determine if an effective address can be referenced with a kernel-key. Note that read/write attributes are not maintained at a page granularity. This service only checks if the kernel-key assigned to an effective address matches the *kkey* value.

## Execution Environment

The `raschk_eaddr_kkey` kernel service can be called from the process or interrupt environment.

## Return Values

Item	Description
0	Successful.
EFAULT_RASCHK_EADDR_KKEY	Operation cannot be performed because a page in or page fault was not allowed.
EINVAL_RASCHK_EADDR_KKEY	The address to validate was determined to be invalid.
EINVAL_RASCHK_EADDR_KKEY_PROT	The address failed the protection check.

### Related reference:

“`raschk_eaddr_hkeyset` Kernel Service” on page 435

## raschk\_stktrace Kernel Service

### Purpose

Generates a runtime compact stack trace for only call chain addresses.

### Syntax

```
#include <sys/raschk.h>
```

```
kernno_t rashchk_stktrace (trcbufsz, flags, trcbuf)
size_t trcbufsz;
long flags;
void * trcbuf;
```

### Parameters

**Item**  
*trcbuflsz*  
*flags*

**Description**

Size of the stack trace buffer the caller allocated.

The following flags are defined:

**RAS\_STK\_DO\_CURMST**

If this flag bit value is set, this service will not look at the previous MST to get the stack trace. The stack trace is obtained only for the current context.

**RAS\_STK\_DO\_PREVMST**

If this flag bit value is set, this service will skip the current MST and start getting the stack trace from the previous MST.

**RAS\_STK\_DO\_ONEMST**

This flag bit value can be combined with the above bit values to get stack trace for that MST.

**RAS\_STK\_GET\_SYMBOLS**

If this flag bit value is set, then all the call chain addresses are translated into a stream of bytes containing symbols with offset (null terminated) and placed in the caller's buffer.

**RAS\_STK\_DO\_CURRWA**

If this flag bit value is set, this service will use the RWA (recovery work area) associated with the current MST to begin the trace back.

**Note:**

The **RAS\_STK\_DO\_PREVMST**, **RAS\_STK\_DO\_CURMST**, and **RAS\_STK\_DO\_CURRWA** flags are mutually exclusive. Specifying the **RAS\_STK\_DO\_ONEMST** flag without specifying the **RAS\_STK\_DO\_PREVMST** flag is equivalent to specifying the **RAS\_STK\_DO\_CURMST** flag.

If the **RAS\_STK\_GET\_SYMBOLS** flag is not set, the end of the stack trace is indicated by an entry containing 0. A value of -2 in *trcbuf* indicates the start of a new **mst** trace if any. Also, the stack trace will stop once we reach the system call boundary as we are interested only in kernel stack trace and we can only validate kernel stack addresses.

If the **RAS\_STK\_GET\_SYMBOLS** flag is set, the output buffer will contain a null-terminated string with the symbolic representation of the stack trace. A call to **raschk\_addr2sym()** is performed for each entry in the stack trace and the resulting strings are concatenated in the output buffer, and separated by '\n' characters. Special values in the stack trace will be translated to appropriate strings.

*trcbuf*

Pointer to the buffer that the caller allocated to get stack trace.

**Note:** Ensure that *trcbuf* is pinned when called disabled.

## Description

This kernel service can be used to generate a runtime compact stack trace. The algorithm is performed for:

- All MSTs starting from the current MST (default, and none of **RAS\_STK\_DO\_CURMST**, **RAS\_STK\_DO\_PREVMST**, **RAS\_STK\_DO\_CURRWA**, nor **RAS\_STK\_DO\_ONEMST** flag bits specified.)
- Only for the current MST (**RAS\_STK\_DO\_CURMST** bit flag is set)
- All the MSTs starting from previous MST (**RAS\_STK\_DO\_PREVMST** bit flag is set)
- Only for the previous MST (**RAS\_STK\_DO\_PREVMST** and **RAS\_STK\_DO\_ONEMST** bits are set)
- For the current MST recovery work area (RWA) context and previous MSTs. (**RAS\_STK\_DO\_CURRWA** flag bit is set.)
- Only for the current MST recovery work area (RWA) context. (**RAS\_STK\_DO\_CURRWA** and **RAS\_STK\_DO\_ONEMST** flag bits are set.)
- Getting all the symbols plus offset corresponding to the call addresses obtained in *trcbuf* and replacing *trcbuf* with symbol information in a string format. (**RAS\_STK\_GET\_SYMBOLS** bit flag is set)

## Execution Environment

The `raschk_stktrace` kernel service can be called from either the process or interrupt environment.

## Return Values

Item	Description
0	Successful
kernno	Unsuccessful

## raw\_input Kernel Service

### Purpose

Builds a `raw_header` structure for a packet and sends both to the raw protocol handler.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
void raw_input (m0, proto, src, dst)
struct mbuf * m0;
struct sockproto * proto;
struct sockaddr * src;
struct sockaddr * dst;
```

### Parameters

Item	Description
<i>m0</i>	Specifies the address of an <code>mbuf</code> structure containing input data.
<i>proto</i>	Specifies the protocol definition of data.
<i>src</i>	Identifies the <code>sockaddr</code> structure indicating where data is from.
<i>dst</i>	Identifies the <code>sockaddr</code> structure indicating the destination of the data.

### Description

The `raw_input` kernel service accepts an input packet, builds a `raw_header` structure (as defined in the `/usr/include/net/raw_cb.h` file), and passes both on to the raw protocol input handler.

## Execution Environment

The `raw_input` kernel service can be called from either the process or interrupt environment.

## Return Values

The `raw_input` service has no return values.

### Related information:

Network Kernel Services

## raw\_usrreq Kernel Service

### Purpose

Implements user requests for raw protocols.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
void raw_usrreq (so, req, m, nam, control)
struct socket * so;
int req;
struct mbuf * m;
struct mbuf * nam;
struct mbuf * control;
```

## Parameters

Item	Description
<i>so</i>	Identifies the address of a raw socket.
<i>req</i>	Specifies the request command.
<i>m</i>	Specifies the address of an <b>mbuf</b> structure containing data.
<i>nam</i>	Specifies the address of an <b>mbuf</b> structure containing the <b>sockaddr</b> structure.
<i>control</i>	This parameter should be set to a null value.

## Description

The **raw\_usrreq** kernel service implements user requests for the raw protocol.

The **raw\_usrreq** service supports the following commands:

Command	Description
<b>PRU_ABORT</b>	Aborts (fast DISCONNECT, DETACH).
<b>PRU_ACCEPT</b>	Accepts connection from peer.
<b>PRU_ATTACH</b>	Attaches protocol to up.
<b>PRU_BIND</b>	Binds socket to address.
<b>PRU_CONNECT</b>	Establishes connection to peer.
<b>PRU_CONNECT2</b>	Connects two sockets.
<b>PRU_CONTROL</b>	Controls operations on protocol.
<b>PRU_DETACH</b>	Detaches protocol from up.
<b>PRU_DISCONNECT</b>	Disconnects from peer.
<b>PRU_LISTEN</b>	Listens for connection.
<b>PRU_PEERADDR</b>	Fetches peer's address.
<b>PRU_RCVD</b>	Have taken data; more room now.
<b>PRU_RCVOOB</b>	Retrieves out of band data.
<b>PRU_SEND</b>	Sends this data.
<b>PRU_SENDOOB</b>	Sends out of band data.
<b>PRU_SENSE</b>	Returns status into m.
<b>PRU_SOCKADDR</b>	Fetches socket's address.
<b>PRU_SHUTDOWN</b>	Will not send any more data.

Any unrecognized command causes the **panic** kernel service to be called.

## Execution Environment

The **raw\_usrreq** kernel service can be called from either the process or interrupt environment.

## Return Values



Item	Description
EOPNOTSUPP	Indicates an unsupported command.
EINVAL	Indicates a parameter error.
EACCES	Indicates insufficient authority to support the PRU_ATTACH command.
ENOTCONN	Indicates an attempt to detach when not attached.
EISCONN	Indicates that the caller tried to connect while already connected.

#### Related reference:

“panic Kernel Service” on page 403

#### Related information:

Network Kernel Services

## reconfig\_register, reconfig\_register\_ext, reconfig\_unregister, or reconfig\_complete, reconfig\_register\_list Kernel Service

### Purpose

Register and unregister reconfiguration handlers.

### Syntax

```
#include <sys/dr.h>
```

```
int reconfig_register (handler, actions,
                      h_arg, h_token, name)
int (*handler)(void *event, void *h_arg, int req,
               void *resource_info);
int actions;
void *h_arg;
ulong *h_token;
char *name;
```

```
int reconfig_register_ext (handler, actions, h_arg, h_token, name)
int (*handler)(void *event, void *h_arg, unsigned long long req,
               void *resource_info);
unsigned long long actions;
void *h_arg;
ulong *h_token;
char *name;
```

```
int reconfig_unregister (h_token)
ulong h_token;
```

```
void reconfig_complete (event, rc)
void *event;
int rc;
```

```
int reconfig_register_list (handler, event_list, list_size, h_arg, h_token, name)
int (*handler)(void *event, void *h_arg, dr_kevent_t event_in_prog,
               void *resource_info);
dr_kevent_t event_list[];
size_t list_size;
void *h_arg;
ulong *h_token;
char *name;
```

### Description

The `reconfig_register`, `reconfig_register_ext`, `reconfig_register_list` and `reconfig_unregister` kernel services register and unregister reconfiguration handlers, which are invoked by the kernel both before and after DLPAR operations depending on the set of events specified by the kernel extension when registering.

Starting with AIX 6.1 with 6100-02, all future kernel extensions use the **reconfig\_register\_list** kernel service when registering for DLPAR operations. The **reconfig\_register\_list** kernel service supports previous and new DLPAR operations. The **reconfig\_register** or **reconfig\_register\_ext** kernel services will no longer support all future DLPAR operations.

The **reconfig\_complete** kernel service is used to indicate that the request has completed. If a kernel extension expects that the operation is likely to take a long time (several seconds), the handler must return **DR\_WAIT** to the caller, but proceed with the request asynchronously. In this case, the handler must indicate that it has completed the request by invoking the **reconfig\_complete** kernel service.

## Parameters

Item	Description
<i>actions</i>	Allows the kernel extension to specify which of the following events require notification: <ul style="list-style-type: none"> <li>• <b>DR_PMIG_CHECK</b></li> <li>• <b>DR_PMIG_PRE</b></li> <li>• <b>DR_PMIG_POST</b></li> <li>• <b>DR_PMIG_POST_ERROR</b></li> <li>• <b>DR_CAP_ADD_CHECK</b></li> <li>• <b>DR_CAP_ADD_PRE</b></li> <li>• <b>DR_CAP_ADD_POST</b></li> <li>• <b>DR_CAP_ADD_POST_ERROR</b></li> <li>• <b>DR_CAP_REMOVE_CHECK</b></li> <li>• <b>DR_CAP_REMOVE_PRE</b></li> <li>• <b>DR_CAP_REMOVE_POST</b></li> <li>• <b>DR_CAP_REMOVE_POST_ERROR</b></li> <li>• <b>DR_CPU_ADD_CHECK</b></li> <li>• <b>DR_CPU_ADD_PRE</b></li> <li>• <b>DR_CPU_ADD_POST</b></li> <li>• <b>DR_CPU_ADD_POST_ERROR</b></li> <li>• <b>DR_CPU_REMOVE_CHECK</b></li> <li>• <b>DR_CPU_REMOVE_PRE</b></li> <li>• <b>DR_CPU_REMOVE_POST</b></li> <li>• <b>DR_CPU_REMOVE_POST_ERROR</b></li> <li>• <b>DR_MEM_ADD_CHECK</b></li> <li>• <b>DR_MEM_ADD_OP_POST</b></li> <li>• <b>DR_MEM_ADD_PRE</b></li> <li>• <b>DR_MEM_ADD_POST</b></li> <li>• <b>DR_MEM_ADD_POST_ERROR</b></li> <li>• <b>DR_MEM_REMOVE_CHECK</b></li> <li>• <b>DR_MEM_REMOVE_OP_POST</b></li> <li>• <b>DR_MEM_REMOVE_OP_PRE</b></li> <li>• <b>DR_MEM_REMOVE_PRE</b></li> <li>• <b>DR_MEM_REMOVE_POST</b></li> <li>• <b>DR_MEM_REMOVE_POST_ERROR</b></li> </ul>
<i>event</i>	Passed to the handler and intended to be used only when calling the <b>reconfig_complete</b> kernel service.
<i>event_list</i>	Specifies which events require notification. For the supported values, see the <b>dr.h</b> file.
<i>handler</i>	Specifies the kernel extension function to be invoked.
<i>h_arg</i>	Specified by the kernel extension, remembered by the kernel along with the function descriptor for the handler, and passed to the handler when it is invoked. It is not used directly by the kernel, but is intended to support kernel extensions that manage multiple adapter instances. This parameter points to an adapter control block.
<i>h_token</i>	An output parameter that is used when unregistering the handler.
<i>list_size</i>	Specifies the memory size of the <b>event_list</b> array.

**Item***name**rc**resource\_info***Description**

Provided for information purposes and may be included within an error log entry, if the driver returns an error. It is provided by the kernel extension and must be limited to 15 ASCII characters.

Can be set to **DR\_FAIL** or **DR\_SUCCESS**.

Identifies the resource specific information for the current DLPAR request. If the request is cpu based, the *resource\_info* data is provided through a **dri\_cpu** structure. Otherwise a **dri\_mem** structure is used. On a Micro-Partitioning partition, if the request is CPU-capacity based, the *resource\_info* data is provided through a **dri\_cpu\_capacity** structure, which has the following format. The kernel extensions are not notified of changes in variable capacity weight in an uncapped Micro-Partitioning environment.

```
*/
struct dri_cpu_capacity {
    uint64_t ent_capacity; /* partition current entitled capacity*/
    int delta_ent_cap; /* delta capacity added/removed*/
    int status; /* capacity update constrained or not */
};

/*
 * dri_cpu_capacity.status flags.
 */
#define CAP_UPDATE_SUCCESS 0x0
#define CAP_UPDATE_CONSTRAINED 0x1
```

**Note:** The capacity update is constrained by the Hypervisor.

If the request is memory capacity based, the *resource\_info* data is provided through a **dri\_mem\_capacity** structure, which has the following format:

```
struct dri_mem_capacity {
    size64_t mem_capacity; /* partition current entitled capacity*/
    ssize64_t delta_mem_capacity;
    uint flags;
    int status; /* capacity update constrained or not */
    uchar reserved[7];
};

/*
 * dri_mem_capacity.status flags.
 */
#define CAP_UPDATE_SUCCESS 0x0
#define CAP_UPDATE_CONSTRAINED 0x1
```

Item	Description
<i>req</i>	Indicates the following DLPAR operation to be performed by the handler: <ul style="list-style-type: none"> <li>• DR_PMIG_CHECK</li> <li>• DR_PMIG_PRE</li> <li>• DR_PMIG_POST</li> <li>• DR_PMIG_POST_ERROR</li> <li>• DR_CAP_ADD_CHECK</li> <li>• DR_CAP_ADD_PRE</li> <li>• DR_CAP_ADD_POST</li> <li>• DR_CAP_ADD_POST_ERROR</li> <li>• DR_CAP_REMOVE_CHECK</li> <li>• DR_CAP_REMOVE_PRE</li> <li>• DR_CAP_REMOVE_POST</li> <li>• DR_CAP_REMOVE_POST_ERROR</li> <li>• DR_CPU_ADD_CHECK</li> <li>• DR_CPU_ADD_PRE</li> <li>• DR_CPU_ADD_POST</li> <li>• DR_CPU_ADD_POST_ERROR</li> <li>• DR_CPU_REMOVE_CHECK</li> <li>• DR_CPU_REMOVE_PRE</li> <li>• DR_CPU_REMOVE_POST</li> <li>• DR_CPU_REMOVE_POST_ERROR</li> <li>• DR_MEM_ADD_CHECK</li> <li>• DR_MEM_ADD_OP_POST</li> <li>• DR_MEM_ADD_PRE</li> <li>• DR_MEM_ADD_POST</li> <li>• DR_MEM_ADD_POST_ERROR</li> <li>• DR_MEM_REMOVE_CHECK</li> <li>• DR_MEM_REMOVE_OP_POST</li> <li>• DR_MEM_REMOVE_OP_PRE</li> <li>• DR_MEM_REMOVE_PRE</li> <li>• DR_MEM_REMOVE_POST</li> <li>• DR_MEM_REMOVE_POST_ERROR</li> </ul>

### List of `dr_kevent_t` events

The following events are used with the `reconfig_register_list()` call for the `event_list` array:

- DR\_KEVENT\_CPU\_ADD\_CHECK
- DR\_KEVENT\_CPU\_ADD\_PRE
- DR\_KEVENT\_CPU\_ADD\_POST
- DR\_KEVENT\_CPU\_ADD\_POST\_ERROR
- DR\_KEVENT\_CPU\_RM\_CHECK
- DR\_KEVENT\_CPU\_RM\_PRE
- DR\_KEVENT\_CPU\_RM\_POST
- DR\_KEVENT\_CPU\_RM\_POST\_ERROR
- DR\_KEVENT\_MEM\_ADD\_CHECK
- DR\_KEVENT\_MEM\_ADD\_PRE
- DR\_KEVENT\_MEM\_ADD\_POST
- DR\_KEVENT\_MEM\_ADD\_POST\_ERROR

- DR\_KEVENT\_MEM\_RM\_CHECK
- DR\_KEVENT\_MEM\_RM\_PRE
- DR\_KEVENT\_MEM\_RM\_POST
- DR\_KEVENT\_MEM\_RM\_POST\_ERROR
- DR\_KEVENT\_MEM\_ADD\_RES
- DR\_KEVENT\_MEM\_RM\_RES
- DR\_KEVENT\_CPU\_CAP\_ADD\_CHECK
- DR\_KEVENT\_CPU\_CAP\_ADD\_PRE
- DR\_KEVENT\_CPU\_CAP\_ADD\_POST
- DR\_KEVENT\_CPU\_CAP\_ADD\_POST\_ERROR
- DR\_KEVENT\_CPU\_CAP\_RM\_CHECK
- DR\_KEVENT\_CPU\_CAP\_RM\_PRE
- DR\_KEVENT\_CPU\_CAP\_RM\_POST
- DR\_KEVENT\_CPU\_CAP\_RM\_POST\_ERROR
- DR\_KEVENT\_MEM\_RM\_OP\_PRE
- DR\_KEVENT\_MEM\_RM\_OP\_POST
- DR\_KEVENT\_MEM\_ADD\_OP\_POST
- DR\_KEVENT\_PMIG\_CHECK
- DR\_KEVENT\_PMIG\_PRE
- DR\_KEVENT\_PMIG\_POST
- DR\_KEVENT\_PMIG\_POST\_ERROR
- DR\_KEVENT\_PMIG\_POST\_INTERNAL
- DR\_KEVENT\_WMIG\_CHECK
- DR\_KEVENT\_WMIG\_PRE
- DR\_KEVENT\_WMIG\_POST
- DR\_KEVENT\_WMIG\_POST\_ERROR
- DR\_KEVENT\_WMIG\_CHECKPOINT\_CHECK
- DR\_KEVENT\_WMIG\_CHECKPOINT\_PRE
- DR\_KEVENT\_WMIG\_CHECKPOINT\_DOIT
- DR\_KEVENT\_WMIG\_CHECKPOINT\_ERROR
- DR\_KEVENT\_WMIG\_CHECKPOINT\_POST
- DR\_KEVENT\_WMIG\_CHECKPOINT\_POST\_ERROR
- DR\_KEVENT\_WMIG\_RESTART\_CHECK
- DR\_KEVENT\_WMIG\_RESTART\_PRE
- DR\_KEVENT\_WMIG\_RESTART\_DOIT
- DR\_KEVENT\_WMIG\_RESTART\_ERROR
- DR\_KEVENT\_WMIG\_RESTART\_POST
- DR\_KEVENT\_WMIG\_RESTART\_POST\_ERROR
- DR\_KEVENT\_MEM\_CAP\_ADD\_CHECK
- DR\_KEVENT\_MEM\_CAP\_ADD\_PRE
- DR\_KEVENT\_MEM\_CAP\_ADD\_POST
- DR\_KEVENT\_MEM\_CAP\_ADD\_POST\_ERROR
- DR\_KEVENT\_MEM\_CAP\_RM\_CHECK
- DR\_KEVENT\_MEM\_CAP\_RM\_PRE
- DR\_KEVENT\_MEM\_CAP\_RM\_POST

- DR\_KEVENT\_MEM\_CAP\_RM\_POST\_ERROR
- DR\_KEVENT\_MEM\_CAP\_WGT\_ADD\_CHECK
- DR\_KEVENT\_MEM\_CAP\_WGT\_ADD\_PRE
- DR\_KEVENT\_MEM\_CAP\_WGT\_ADD\_POST
- DR\_KEVENT\_MEM\_CAP\_WGT\_ADD\_POST\_ERROR
- DR\_KEVENT\_MEM\_CAP\_WGT\_RM\_CHECK
- DR\_KEVENT\_MEM\_CAP\_WGT\_RM\_PRE
- DR\_KEVENT\_MEM\_CAP\_WGT\_RM\_POST
- DR\_KEVENT\_MEM\_CAP\_WGT\_RM\_POST\_ERROR
- DR\_KEVENT\_TOPOLOGY\_PRE
- DR\_KEVENT\_TOPOLOGY\_POST
- DR\_KEVENT\_AME\_FACTOR\_CHECK
- DR\_KEVENT\_AME\_FACTOR\_PRE
- DR\_KEVENT\_AME\_FACTOR\_POST
- DR\_KEVENT\_AME\_FACTOR\_POST\_ERROR

## Return Values

Upon successful completion, the `reconfig_register`, `reconfig_register_ext` and `reconfig_unregister` kernel services return zero. If unsuccessful, the appropriate `errno` value is returned.

## Execution Environment

The `reconfig_register`, `reconfig_register_ext`, `reconfig_unregister`, and `handler` interfaces are invoked in the process environment only.

The `reconfig_complete` kernel service may be invoked in the process or interrupt environment.

### Related information:

Making Kernel Extensions DLPAR-Aware

## refmon Kernel Service

### Purpose

Performs various access checks such as privileges, authorizations, discretionary access control checks and so on.

### Syntax

```
#include <refmon.h>
```

```
int refmon (crp, action, flags, nargs, args[])
cred_t *crp;
rfm_action_t action;
uint_t flags;
int nargs;
void *args[];
```

### Parameters

Item	Description
<i>crp</i>	Specifies the caller's (subject) credentials; If NULL, then current process credentials are referenced.
<i>action</i>	Specifies the type of access check that needs to be carried out.
<i>flags</i>	Enables auditing of this event. You can only set this parameter to the value of REFMON_AUDIT.
<i>nargs</i>	Specifies the number of arguments in the <i>args</i> parameter.
<i>args</i>	Specifies an array of void pointers used as input to the <b>refmon</b> kernel service based on the <i>action</i> parameter.

## Description

The **refmon** kernel service provides an interface to perform various access checks. You can call the **refmon** kernel service to determine access to system resources. Most of the actions that are passed to the **refmon** kernel service check for specific privileges. Many of the system calls and kernel services call the **refmon** kernel service to check whether you are authorized or privileged to use such functions. The *action* parameter determines which type of checks needs to be performed. The **sys/refmon.h** header file contains a complete list of these actions and their corresponding description.

## Execution Environment

The **refmon** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Success.
EINVAL	The <i>action</i> parameter is not valid or a value that is not allowed is passed in for an action.
EPERM	The caller does not have permission to perform the intended action.

### Related information:

Security Kernel Services

## register\_HA\_handler Kernel Service Purpose

Registers a High Availability Event Handler with the Kernel.

## Syntax

```
#include <sys/high_avail.h>
```

```
int register_HA_handler (ha_handler)
ha_handler_ext_t * ha_handler;
```

## Parameter

Item	Description
<i>ha_handler</i>	Specifies a pointer to a structure of the type <b>ha_handler_ext_t</b> as defined in <b>/usr/include/sys/high_avail.h</b> .

## Description

The **register\_HA\_handler** kernel registers the **High Availability Event Handler (HAEH)** function to those kernel extensions that need to be made aware of high availability events such as processor deallocation. This function is called by the kernel, at base level, when a high availability event is initiated, due to some hardware fault.

The **ha\_handler\_ext\_t** structure has 3 fields:

Field	Description
<i>_fun</i>	Contains a pointer to the high availability event handler function.
<i>_data</i>	Contains a user defined value which will be passed as an argument by the kernel when calling the function.
<i>_name</i>	Component name

When a high availability event is initiated, the kernel calls *\_fun()* at base level (that is, process environment) with 2 parameters:

- The first is the data the user passed in the *\_data* field at registration time.
- The second is a pointer to a **haeh\_event\_t** structure defined in */usr/include/sys/high\_avail.h*.

The fields of interest in this structure are:

Field	Description
<i>_magic</i>	Identifies the event type. The only possible value is <b>HA_CPU_FAIL</b> .
<i>dealloc_cpu</i>	The logical number of the CPU being deallocated.

The high availability even handler, in addition to user specific functions, must unbind its threads bound to *dealloc\_cpu* and stop the timer request blocks (TRB) started by those bound threads when applicable.

The high availability event handler must return one of the following values:

Value	Description
<b>HA_ACCEPTED</b>	The user processing of the event has succeeded.
<b>HA_REFUSED</b>	The user processing of the event was not successful.

Any return value different from **HA\_ACCEPTED** causes the kernel to abort the processing of the event. In the case of a processor failure, the processor deallocation is aborted. In this case, a **CPU\_DEALLOC\_ABORTED** error log entry is created, and the value passed in the *\_name* field appears in the detailed data area of the error log entry.

An extension may register the same HAEH *N* times (*N* > 1). Although it is considered as an incorrect behaviour, no error is reported. The given HAEH is invoked *N* times for each HA event. This handler has to be unregistered as many times as it was registered.

Since the kernel calls the HAEH in turn, it is possible for a HAEH to be called multiple times for the same event. The kernel extensions should be ready to deal with this possibility. For example, two kernel extensions **K1** and **K2** have registered HA Handlers. A CPU deallocation is initiated. The HAEH for **K1** gets invoked, does its job and returns **HA\_ACCEPTED**. **K2** gets invoked next and for some reason returns **HA\_REFUSED**. The deallocation is aborted, and an error log entry reports **K2** as the reason for failure. Later, the system administrator unloads **K2** and restarts the deallocation by manually running **ha\_star**. The result is that the HAEH for **K1** gets invoked again with the same parameters.

## Execution Environment

The **register\_HA\_handler** kernel service can be called from the process environment only.

## Return Values



Item	Description
0	Indicates a successful operation.

A non zero value indicates an error.

**Related reference:**

“unregister\_HA\_handler Kernel Service” on page 522

**Related information:**

RAS Kernel Services

## rmalloc Kernel Service

### Purpose

Allocates an area of memory from the **real\_heap** heap.

### Syntax

```
#include <sys/types.h>
caddr_t rmalloc (size, align)
int size
int align
```

### Parameters

Item	Description
<i>size</i>	Specifies the number of bytes to allocate.
<i>align</i>	Specifies alignment characteristics.

### Description

The **rmalloc** kernel service allocates an area of memory from the contiguous real memory heap. This area is the number of bytes in length specified by the *size* parameter and is aligned on the byte boundary specified by the *align* parameter. The *align* parameter is actually the log base 2 of the desired address boundary. For example, an *align* value of 4 requests that the allocated area be aligned on a 16-byte boundary.

The contiguous real memory heap, **real\_heap**, is a heap of contiguous real memory pages located in the low 16MB of real memory. This heap is virtually mapped into the kernel extension's address space. By nature, this heap is implicitly pinned, so no explicit pinning of allocated regions is necessary.

The **real\_heap** heap is useful for devices that require DMA transfers greater than 4K but do not provide a scatter/gather capability. Such a device must be given contiguous bus addresses by its device driver. The device driver should pass the **DMA\_CONTIGUOUS** flag on its **d\_map\_init** call in order to obtain contiguous mappings. On certain platforms it is possible that a **d\_map\_init** call using the **DMA\_CONTIGUOUS** flag could fail. In this case, the device driver can make use of the **real\_heap** heap (using **rmalloc**) to obtain contiguous bus addresses for its device driver. Because the **real\_heap** heap is a limited resource, device drivers should always attempt to use the **DMA\_CONTIGUOUS** flag first.

On unsupported platforms, the **rmalloc** service returns NULL if the requested memory cannot be allocated.

The **rmfree** kernel service should be called to free allocation from a previous **rmalloc** call. The **rmalloc** kernel service can be called from the process environment only.

## Return Values

Upon successful completion, the **rmalloc** kernel service returns the address of the allocated area. A **NULL** pointer is returned if the requested memory cannot be allocated.

**Related reference:**

“rmfree Kernel Service”

## rmfree Kernel Service

### Purpose

Frees memory allocated by the **rmalloc** kernel service.

### Syntax

```
#include <sys/types.h>
```

```
int rmfree ( pointer, size)
caddr_t pointer
int size
```

### Parameters

Item	Description
<i>pointer</i>	Specifies the address of the area in memory to free.
<i>size</i>	Specifies the size of the area in memory to free.

### Description

The **rmfree** kernel service frees the area of memory pointed to by the *pointer* parameter in the contiguous real memory heap. This area of memory must be allocated with the **rmalloc** kernel service, and the *pointer* must be the pointer returned from the corresponding **rmalloc** kernel service call. Also, the *size* must be the same size that was used on the corresponding **rmalloc** call.

Any memory allocated in a prior **rmalloc** call must be explicitly freed with an **rmfree** call. This service can be called from the process environment only.

### Return Values

Item	Description
0	Indicates successful completion.
-1	Indicates one of the following: <ul style="list-style-type: none"><li>• The area was not allocated by the <b>rmalloc</b> kernel service.</li><li>• The heap was not initialized for memory allocation.</li></ul>

**Related reference:**

“rmalloc Kernel Service” on page 449

## rmmmap\_create Kernel Service

### Purpose

Defines an Effective Address [EA] to Real Address [RA] translation region.

### Syntax

```
#include <sys/ioacc.h>
#include <sys/adspace.h>
```

```

int rmmmap_create ( eaddrp, iomp, flags)
void **eaddrp;
struct io_map *iomp;
int flags;

```

## Parameters

Item	Description
<i>eaddr</i>	Required process effective address of the mapping region.
<i>iomp</i>	The bus memory to which the effective address described by the <i>eaddr</i> parameter must correspond. For real memory, the bus id must be set to <b>REALMEM_BID</b> and the bus address must be set to the real memory address. The size field must be at least <b>PAGESIZE</b> , no larger than <b>SEGSIZE</b> , and a multiple of <b>PAGESIZE</b> . The key must be set to <b>IO_MEM_MAP</b> . The flags field is not used.
<i>flags</i>	The flags select page and segment attributes of the translation. Not all page attribute flags are compatible. The valid combinations of page attribute flags follow. <ul style="list-style-type: none"> <li><b>RMMAP_PAGE_W</b> PowerPC "Write Through" page attribute. Write-through mode is not supported, and if this flag is set, <b>EINVAL</b> is reported.</li> <li><b>RMMAP_PAGE_I</b> PowerPC "Cache Inhibited" page attribute. This flag is valid for I/O mappings, but is not allowed for real memory mappings.</li> <li><b>RMMAP_PAGE_M</b> PowerPC "Memory Coherency Required" page attribute. This flag is optional for I/O mappings; however, it is required for memory mappings. The default operating mode for real memory mappings has this bit set.</li> <li><b>RMMAP_PAGE_G</b> PowerPC "Guarded" page attribute. This flag is optional for I/O mappings, and must be 0 for real memory mappings. Although optional for I/O, it is recommended that this flag must be set for I/O mappings. When set, the processor does not make unnecessary (speculative) references to the page. It includes out of order read or write operations and branch fetching. When clear, normal PowerPC speculative execution rules apply.</li> <li><b>RMMAP_RDONLY</b> When set, the page protection bits used in the <b>HTAB</b> does not allow write operations regardless of the setting of the key bit in the associated segment register. Exactly one of <b>RMMAP_RDONLY</b> and <b>RMMAP_RDWR</b> must be specified.</li> <li><b>RMMAP_RDWR</b> When set, the page protection bits used in the <b>HTAB</b> allows read and write operations regardless of the setting of the key bit in the associated segment register. Exactly one of: <b>RMMAP_RDONLY</b>, and <b>RMMAP_RDWR</b> must be specified.</li> <li><b>RMMAP_PRELOAD</b> When set, the protection attributes of this region are entered immediately into the hardware page table. It is very slow initially, but prevents each referenced page in the region from faulting in separately. It is only advisory. This flag is not maintained as an attribute of the map region, it is used only during the current call.</li> <li><b>RMMAP_INHERIT</b> When set, this protection attribute specifies that the translation region created by this <b>rmmmap_create</b> invocation must be inherited on a <b>fork</b> operation, to the child process. This inheritance is achieved with copy-semantics. The child has its own private mapping to the same I/O or real memory address range as the parent.</li> </ul>

## Description

The translation regions that are created with **rmmmap\_create** kernel service are maintained in I/O mapping segments. Any single such segment might translate up to 256 Megabytes of real memory or memory mapped I/O in a single region. The only granularity for which the **rmmmap\_remove** service might be started is a single mapping that is created by a single call to the **rmmmap\_create**.

There are constraints on the size of the mapping and the *flags* parameter, described later, which causes the call to fail regardless of whether adequate effective address space exists.

If **rmmmap\_create** kernel service is called with the effective address of zero, the function attempts to find free space in the process address space. If successful, an I/O mapping segment is created and the effective address (which is passed by reference) is changed to the effective address which is mapped to the first page of the *iomp* memory.

If **rmmmap\_create** kernel service is called with a non-zero effective address, it is taken as the required effective address which must translate to the passed *iomp* memory. This function verifies that the requested range is free. If not, it fails and returns **EINVAL**. If the mapping at the effective address is not contained in a single segment, the function fails and returns **ENOSPC**. Otherwise, the region is allocated and the effective address is not modified. The effective address is mapped to the first page of the *iomp* memory. References outside of the mapped regions but within the same segment are invalid.

The effective address (if provided) and the bus address must be a multiple of **PAGESIZE** or **EINVAL** is returned.

I/O mapping segments are not inherited by child processes after a **fork** subroutine.

I/O mapping segments are not inherited by child processes after a **fork** subroutine, except when **RMMAP\_INHERIT** is specified. These segments are deleted by **exec**, **exit**, or **rmmmap\_remove** of the last range in a segment.

Only certain combinations of flags are permitted, depending on the type of memory that is mapped. For real memory mappings, **RMMAP\_PAGE\_M** is required while **RMMAP\_PAGE\_W**, **RMMAP\_PAGE\_I**, and **RMMAP\_PAGE\_G** are not allowed. For I/O mappings, it is valid to specify only **RMMAP\_PAGE\_M**, with no other page attribute flags. It is also valid to specify **RMMAP\_PAGE\_I** and optionally, either or both of **RMMAP\_PAGE\_M**, and **RMMAP\_PAGE\_G**. **RMMAP\_PAGE\_W** is never allowed.

The real address range that is described by the *iomp* parameter must be unique within this I/O mapping segment.

## Execution Environment

The **rmmmap\_create** kernel service can be called only from the process environment.

## Return Values

On successful completion, **rmmmap\_create** kernel service returns zero and modifies the effective address to the value at which the newly created mapping region was attached to the process address space. Otherwise, it returns one of following errors:

Item	Description
<b>EINVAL</b>	Some type of parameter error occurred. These parameters include, but are not limited to, size errors and mutually exclusive flag selections.
<b>ENOMEM</b>	The operating system cannot allocate the necessary data structures to represent the mapping.
<b>ENOSPC</b>	Effective address space exhausted in the region indicated by <i>eaddr</i> .
<b>EPERM</b>	This hardware platform does not implement this service.

## Implementation Specifics

This service only functions on PowerPC microprocessors.

### Related reference:

“**rmmmap\_remove** Kernel Service” on page 454

### Related information:

Memory Kernel Services

## **rmmmap\_getwing Kernel Service**

### **Purpose**

Returns wing information about a particular effective address range within an effective address to real address translation region.

### **Syntax**

```
#include <sys/adspage.h>
int rmmmap_getwing(eaddr, npages, results)
unsigned long long eaddr;
unsigned int npages;
char* results;
```

### **Parameters**

<b>Item</b>	<b>Description</b>
<i>eaddr</i>	The process effective address of the start of the desired mapping region. This address should point somewhere inside the first page of the range. This address is interpreted as a 64-bit quantity if the current user address space is 64-bits, and is interpreted as a 32-bit (not remapped) quantity if the current user address space is 32-bits.
<i>npages</i>	The number of pages whose wing information is returned, starting from the page indicated by <b>eaddr</b> .
<i>results</i>	This is an array of bytes, where the wing information is returned. The address of this is passed in by the caller, and <b>rmmmap_getwing</b> stores the wing information for each page in the range in each successive byte in this array. The size of this array is indicated by <i>npages</i> as specified by the caller. The caller is responsible for ensuring that the storage allocated for this array is large enough to hold <i>npages</i> bytes.

### **Description**

The wing information corresponding to the input effective address range is returned.

This routine only works for regions previously mapped with an I/O mapping segment as created by **rmmmap\_create**.

**npages** should not be such that the range crosses a segment boundary. If it does, EINVAL is returned.

The wing information is returned in the **results** array. Each element of the **results** array is a character. Each character may be added with the following fields to examine wing information:

**RMMAP\_PAGE\_W**, **RMMAP\_PAGE\_I**, **RMMAP\_PAGE\_M** or **RMMAP\_PAGE\_G**. The array is valid if the return value is 0.

### **Execution Environment**

The **rmmmap\_getwing** kernel service is called from the process environment only.

### **Return Values**

Item	Description
0	Successful completion. Indicates that the <i>results</i> array is valid and should be examined.
EINVAL	An error occurred. Most likely the region was not mapped via <code>rmmap_create</code> previously..
EINVAL	Input range crosses a certain boundary.
EINVAL	The hardware platform does not implement this service.

## Implementation Specifics

This service only functions on PowerPC microprocessors.

### Related reference:

“`rmmap_create` Kernel Service” on page 450

“`rmmap_remove` Kernel Service”

## rmmap\_remove Kernel Service

### Purpose

Destroys an effective address to real address translation region.

### Syntax

```
#include <sys/adspace.h>
int rmmap_remove (eaddrp);
void **eaddrp;
```

### Parameters

Item	Description
<i>eaddrp</i>	Pointer to the process effective address of the desired mapping region.

### Description

Destroys an effective address to real address translation region. If `rmmap_remove` kernel service is called with the effective address within the region of a previously created I/O mapping segment, the region is destroyed. This service must be called from the process level.

### Execution Environment

The `rmmap_remove` kernel service can be called from the process environment only.

### Return Values

Item	Description
0	Indicates a successful operation.
EINVAL	The provided <i>eaddr</i> does not correspond to a valid I/O mapping segment.
EINVAL	This hardware platform does not implement this service.

## Implementation Specifics

This service only functions on PowerPC microprocessors.

### Related reference:

“`rmmap_create` Kernel Service” on page 450

### Related information:

Memory Kernel Services

## **rtalloc Kernel Service**

### **Purpose**

Allocates a route.

### **Syntax**

```
#include <sys/types.h>
#include <sys/errno.h>
#include <net/route.h>
```

```
void rtalloc ( ro)
register struct route *ro;
```

### **Parameter**

Item	Description
<i>ro</i>	Specifies the route.

### **Description**

The **rtalloc** kernel service allocates a route, which consists of a destination address and a reference to a routing entry.

### **Execution Environment**

The **rtalloc** kernel service can be called from either the process or interrupt environment.

### **Return Values**

The **rtalloc** service has no return values.

### **Example**

To allocate a route, invoke the **rtalloc** kernel service as follows:

```
rtalloc(ro);
```

#### **Related information:**

Network Kernel Services

## **rtalloc\_gr Kernel Service**

### **Purpose**

Allocates a route.

### **Syntax**

```
#include <sys/types.h>
#include <sys/errno.h>
#include <net/route.h>
```

```
void rtalloc_gr ( ro, gidlist)
register struct route *ro;
struct gidstruct *gidlist;
```

## Parameter

Item	Description
<i>ro</i>	Specifies the route.
<i>gidlist</i>	Points to the group list.

## Description

The **rtalloc\_gr** kernel service allocates a route, which consists of a destination address and a reference to a routing entry.

A route can be allocated only if its group id restrictions specify that it can be used by a user with the *gidlist* that is passed in.

## Execution Environment

The **rtalloc\_gr** kernel service can be called from either the process or interrupt environment.

## Return Values

The **rtalloc\_gr** service has no return values.

## Example

To allocate a route, invoke the **rtalloc\_gr** kernel service as follows:

```
rtalloc_gr (ro, gidlist);
```

### Related reference:

“rtalloc Kernel Service” on page 455

### Related information:

Network Kernel Services

## rtfree Kernel Service

### Purpose

Frees the routing table entry.

### Syntax

```
#include <sys/types.h>  
#include <sys/errno.h>  
#include <net/route.h>
```

```
int rtfree ( rt)  
register struct rtentry *rt;
```

## Parameter



Item	Description
<i>rt</i>	Specifies the routing table entry.

## Description

The **rtfree** kernel service frees the entry it is passed from the routing table. If the route does not exist, the **panic** service is called. Otherwise, the **rtfree** service frees the **mbuf** structure that contains the route and decrements the routing reference counters.

## Execution Environment

The **rtfree** kernel service can be called from either the process or interrupt environment.

## Return Values

The **rtfree** kernel service has no return values.

## Example

To free a routing table entry, invoke the **rtfree** kernel service as follows:

```
rtfree(rt);
```

### Related reference:

“panic Kernel Service” on page 403

### Related information:

Network Kernel Services

## rtinit Kernel Service

### Purpose

Sets up a routing table entry typically for a network interface.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/socket.h>
#include <net/route.h>
```

```
int rtinit (ifa, cmd, flags)
```

```
struct ifaddr * ifa;
```

```
int cmd, flags;
```

### Parameters

Item	Description
<i>ifa</i>	Specifies the address of an <b>ifaddr</b> structure containing destination address, interface address, and netmask.
<i>cmd</i>	Specifies a request to add or delete route entry.
<i>flags</i>	Identifies routing flags, as defined in the <code>/usr/include/net/route.h</code> file.

## Description

The **rtinit** kernel service creates a routing table entry for an interface. It builds an **rtentry** structure using the values in the *ifa* and *flags* parameters.

The **rtinit** service then calls the **rtrequest** kernel service and passes the *cmd* parameter and the **rtentry** structure to process the request. The *cmd* parameter contains either the value **RTM\_ADD** (a request to add the route entry) or the value **RTM\_DELETE** (delete the route entry). Valid routing flags to set are defined in the `/usr/include/route.h` file.

## Execution Environment

The **rtinit** kernel service can be called from either the process or interrupt environment.

## Return Values

The **rtinit** kernel service returns values from the **rtrequest** kernel service.

## Example

To set up a routing table entry, invoke the **rtinit** kernel service as follows:

```
rtinit(ifa, RMT_ADD, flags ( RTF_DYNAMIC);
```

### Related reference:

“rtrequest Kernel Service” on page 459

### Related information:

Network Kernel Services

## rtredirect Kernel Service

### Purpose

Forces a routing table entry with the specified destination to go through a given gateway.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/mbuf.h>
#include <net/route.h>
```

```
rtredirect ( dst, gateway, netmask, flags, src, rtp)
struct sockaddr *dst, *gateway, *netmask, *src;
int flags;
struct rtentry **rtp;
```

### Parameters

Item	Description
<i>dst</i>	Specifies the destination address.
<i>gateway</i>	Specifies the gateway address.
<i>netmask</i>	Specifies the network mask for the route.
<i>flags</i>	Indicates routing flags as defined in the <code>/usr/include/net/route.h</code> file.
<i>src</i>	Identifies the source of the redirect request.
<i>rtp</i>	Indicates the address of a pointer to a <b>rtentry</b> structure. Used to return a constructed route.

### Description

The **rtredirect** kernel service forces a routing table entry for a specified destination to go through the given gateway. Typically, the **rtredirect** service is called as a result of a routing redirect message from the network layer. The *dst*, *gateway*, and *flags* parameters are passed to the **rtrequest** kernel service to process the request.

## Execution Environment

The `rtredirect` kernel service can be called from either the process or interrupt environment.

## Return Values

Item	Description
0	Indicates a successful operation.

If a bad redirect request is received, the routing statistics counter for bad redirects is incremented.

## Example

To force a routing table entry with the specified destination to go through the given gateway, invoke the `rtredirect` kernel service:

```
rtredirect(dst, gateway, netmask, flags, src, rtp);
```

### Related reference:

“rtinit Kernel Service” on page 457

### Related information:

Network Kernel Services

## rtrequest Kernel Service

### Purpose

Carries out a request to change the routing table.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/mbuf.h>
#include <net/if.h>
#include <net/af.h>
#include <net/route.h>
```

```
int rtrequest ( req, dst, gateway, netmask, flags, ret_nrt)
int req;
struct sockaddr *dst, *gateway, *netmask;
int flags;
struct rtable **ret_nrt;
```

### Parameters

Item	Description
<i>req</i>	Specifies a request to add or delete a route.
<i>dst</i>	Specifies the destination part of the route.
<i>gateway</i>	Specifies the gateway part of the route.
<i>netmask</i>	Specifies the network mask to apply to the route.
<i>flags</i>	Identifies routing flags, as defined in the <code>/usr/include/net/route.h</code> file.
<i>ret_nrt</i>	Specifies to return the resultant route.

### Description

The `rtrequest` kernel service carries out a request to change the routing table. Interfaces call the `rtrequest` service at boot time to make their local routes known for routing table ioctl operations. Interfaces also call the `rtrequest` service as the result of routing redirects. The request is either to add (if the *req* parameter

has a value of `RMT_ADD`) or delete (the `req` parameter is a value of `RMT_DELETE`) the route.

## Execution Environment

The `rtrequest` kernel service can be called from either the process or interrupt environment.

## Return Values

Item	Description
0	Indicates a successful operation.
ESRCH	Indicates that the route was not there to delete.
EEXIST	Indicates that the entry the <code>rtrequest</code> service tried to add already exists.
ENETUNREACH	Indicates that the <code>rtrequest</code> service cannot find the interface for the route.
ENOBUFS	Indicates that the <code>rtrequest</code> service cannot get an <code>mbuf</code> structure to add an entry.

## Example

To carry out a request to change the routing table, invoke the `rtrequest` kernel service as follows:

```
rtrequest(RTM_ADD, dst, gateway, netmask, flags, &rtp);
```

### Related reference:

“rtinit Kernel Service” on page 457

### Related information:

Network Kernel Services

## `rtrequest_gr` Kernel Service

### Purpose

Carries out a request to change the routing table.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/mbuf.h>
#include <net/if.h>
#include <net/af.h>
#include <net/route.h>
```

```
int rtrequest_gr ( req, dst, gateway, netmask, flags, ret_nrt, rt_parm)
int req;
struct sockaddr *dst, *
gateway
, *netmask;
int flags;
struct rtable **
ret_nrt;
struct rtreq_parm *
rt_parm;
```

### Parameters

Item	Description
<i>req</i>	Specifies a request to add or delete a route.
<i>dst</i>	Specifies the destination part of the route.
<i>gateway</i>	Specifies the gateway part of the route.
<i>netmask</i>	Specifies the network mask to apply to the route.
<i>flags</i>	Identifies routing flags, as defined in the <code>/usr/include/net/route.h</code> file.
<i>ret_nrt</i>	Specifies to return the resultant route.
<i>rt_parm</i>	Points to the <code>rtreq_parm</code> structure. The <code>/usr/include/net/radix.h</code> file contains the <code>rtreq_parm</code> structure. Through this structure, the route attributes like group list, policy, weight, WPAR ID, interface can be specified.

## Description

The `rtrequest_gr` kernel service carries out a request to change the routing table. Interfaces call the `rtrequest_gr` service at boot time to make their local routes known for routing table ioctl operations. Interfaces also call the `rtrequest_gr` service as the result of routing redirects. The request is either to add (if the *req* parameter has a value of `RMT_ADD`) or delete (the *req* parameter is a value of `RMT_DELETE`) the route.

## Execution Environment

The `rtrequest_gr` kernel service can be called from either the process or interrupt environment.

## Return Values

Item	Description
0	Indicates a successful operation.
ESRCH	Indicates that the route was not there to delete.
EEXIST	Indicates that the entry the <code>rtrequest_gr</code> service tried to add already exists.
ENETUNREACH	Indicates that the <code>rtrequest_gr</code> service cannot find the interface for the route.
ENOBUFS	Indicates that the <code>rtrequest_gr</code> service cannot get an <code>mbuf</code> structure to add an entry.

## Example

To carry out a request to change the routing table, invoke the `rtrequest_gr` kernel service as follows:

```
rtrequest_gr(RTM_ADD, dst, gateway, netmask, flags, &rtp, &rtreq);
```

### Related reference:

“rtinit Kernel Service” on page 457

“rtrequest Kernel Service” on page 459

### Related information:

Network Kernel Services

## rusage\_incr Kernel Service

### Purpose

Increments a field of the `rusage` structure.

### Syntax

```
#include <sys/encap.h>
```

```
void rusage_incr ( field, amount)
int field;
int amount;
```

## Parameters

Item	Description
<i>field</i>	Specifies the field to increment. It must have one of the following values: <b>RUSAGE_INBLOCK</b> Denotes the <code>ru_inblock</code> field. This field specifies the number of times the file system performed input. <b>RUSAGE_OUTBLOCK</b> Denotes the <code>ru_outblock</code> field. This field specifies the number of times the file system performed output. <b>RUSAGE_MSGRCV</b> Denotes the <code>ru_msgrcv</code> field. This field specifies the number of IPC messages received. <b>RUSAGE_MSGSENT</b> Denotes the <code>ru_msgsnd</code> field. This field specifies the number of IPC messages sent.
<i>amount</i>	Specifies the amount to increment to the field.

## Description

The `rusage_incr` kernel service increments the field specified by the *field* parameter of the calling process' `rusage` structure by the amount *amount*.

## Execution Environment

The `rusage_incr` kernel service can be called from the process environment only.

## Return Values

The `rusage_incr` kernel service has no return values.

### Related information:

getrusage subroutine

Process and Exception Management Kernel Services

## S

The following kernel services begin with the with the letter s.

### schednetisr Kernel Service

#### Purpose

Schedules or invokes a network software interrupt service routine.

#### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <net/netisr.h>
```

```
int schednetisr ( anISR )
int anISR;
```

#### Parameter

Item	Description
<i>anISR</i>	Specifies the software interrupt number to issue. Refer to <b>netisr.h</b> for the range of values of <i>anISR</i> that are already in use. Also, other kernel extensions that are not AIX and that use network ISRs currently running on the system can make use of additional values not mentioned in <b>netisr.h</b> .

## Description

The **schednetisr** kernel service schedules or calls a network interrupt service routine. The **add\_netisr** kernel service establishes interrupt service routines. If the service was added with a service level of **NET\_OFF\_LEVEL**, the **schednetisr** kernel service directly calls the interrupt service routine. If the service level was **NET\_KPROC**, a network kernel dispatcher is notified to call the interrupt service routine.

## Execution Environment

The **schednetisr** kernel service can be called from either the process or interrupt environment.

## Return Values

Item	Description
<b>EFAULT</b>	Indicates that a network interrupt service routine does not exist for the specified interrupt number.
<b>EINVAL</b>	Indicates that the <i>anISR</i> parameter is out of range.

### Related reference:

“add\_netisr Kernel Service” on page 10

“del\_netisr Kernel Service” on page 66

### Related information:

Network Kernel Services

## selnotify Kernel Service

### Purpose

Wakes up processes waiting in a **poll** or **select** subroutine or in the **fp\_poll** kernel service.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
void selnotify ( id, subid, rtnevents)
int id;
int subid;
ushort rtnevents;
```

### Parameters

Item	Description
<i>id</i>	Indicates a primary resource identification value. This value along with the subidentifier (specified by the <i>subid</i> parameter) is used by the kernel to notify the appropriate processes of the occurrence of the indicated events. If the resource on which the event has occurred is a device driver, this parameter must be the device major/minor number (that is, a <b>dev_t</b> structure that has been cast to an <b>int</b> ). The kernel has reserved values for the <i>id</i> parameter that do not conflict with possible device major or minor numbers for sockets, message queues, and named pipes.
<i>subid</i>	Helps identify the resource on which the event has occurred for the kernel. For a multiplexed device driver, this is the number of the channel on which the requested events occurred. If the device driver is nonmultiplexed, the <i>subid</i> parameter must be set to 0.
<i>rtnevents</i>	Consists of a set of bits indicating the requested events that have occurred on the specified device or channel. These flags have the same definition as the event flags that were provided by the <i>events</i> parameter on the unsatisfied call to the object's select routine.

## Description

The **selnotify** kernel service should be used by device drivers that support select or poll operations. It is also used by the kernel to support select or poll requests to sockets, named pipes, and message queues.

The **selnotify** kernel service wakes up processes waiting on a **select** or **poll** subroutine. The processes to be awakened are those specifying the given device and one or more of the events that have occurred on the specified device. The **select** and **poll** subroutines allow a process to request information about one or more events on a particular device. If none of the requested events have yet happened, the process is put to sleep and re-awakened later when the events actually happen.

The **selnotify** service should be called whenever a previous call to the device driver's **ddselect** entry point returns and both of the following conditions apply:

- The status of all requested events is false.
- Asynchronous notification of the events is requested.

The **selnotify** service can be called for other than these conditions but performs no operation.

## Sequence of Events for Asynchronous Notification

The device driver must store information about the events requested while in the driver's **ddselect** routine under the following conditions:

- None of the requested events are true (at the time of the call).
- The **POLLSYNC** flag is not set in the *events* parameter.

The **POLLSYNC** flag, when not set, indicates that asynchronous notification is desired. In this case, the **selnotify** service should be called when one or more of the requested events later becomes true for that device and channel.

When the device driver finds that it can satisfy a **select** request, (perhaps due to new input data) and an unsatisfied request for that event is still pending, the **selnotify** service is called with the following items:

- Device major and minor number specified by the *id* parameter
- Channel number specified by the *subid* parameter
- Occurred events specified by the *rtnevents* parameter

These parameters describe the device instance and requested events that have occurred on that device. The notifying device driver then resets its requested-events flags for the events that have occurred for that device and channel. The reset flags thus indicate that those events are no longer requested.

If the *rtnevents* parameter indicated by the call to the **selnotify** service is no longer being waited on, no processes are awakened.

## Execution Environment

The **selnotify** kernel service can be called from either the process or interrupt environment.

## Return Values

The **selnotify** service has no return values.



## Implementation Specifics

The **selnotify** kernel service is part of Base Operating System (BOS) Runtime.

### Related reference:

“ddselect Device Driver Entry Point” on page 633

“fp\_poll Kernel Service” on page 161

“fp\_select Kernel Service” on page 167

“selreg Kernel Service”

### Related information:

poll subroutine

select subroutine

Kernel Extension and Device Driver Management Kernel Services

## selreg Kernel Service

### Purpose

Registers an asynchronous poll or select request with the kernel.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/poll.h>
```

```
int selreg ( corl, dev_id, unique_id, regevents, notify)
int corl;
int dev_id;
int unique_id;
ushort regevents;
void (*notify) ( );
```

### Parameters

Item	Description
<i>corl</i>	The correlator for the poll or select request. The <i>corl</i> parameter is used by the <b>poll</b> and <b>select</b> subroutines to correlate the returned events in a specific select control block with a process' file descriptor or message queue.
<i>dev_id</i>	Primary resource identification value. Along with the <i>unique_id</i> parameter, the <i>dev_id</i> parameter is used to record in the select control block the resource on which the requested poll or select events are expected to occur.
<i>unique_id</i>	Unique resource identification value. Along with the <i>dev_id</i> parameter, the <i>unique_id</i> parameter denotes the resource on which the requested events are expected to occur. For a multiplexed device driver, this parameter specifies the number of the channel on which the requested events are expected to occur. For a nonmultiplexed device driver, this parameter must be set to 0.
<i>regevents</i>	Requested events parameter. The <i>regevents</i> parameter consists of a set of bit flags denoting the events for which notification is being requested. These flags have the same definitions as the event flags provided by the <i>events</i> parameter on the unsatisfied call to the object's <b>select</b> subroutine (see the <b>sys/poll.h</b> file for the definitions). <b>Note:</b> The <b>POLLSYNC</b> bit flag should not be set in this parameter.
<i>notify</i>	Notification routine entry point. This parameter points to a notification routine used for nested poll and select calls.

### Description

The **selreg** kernel service is used by **select** file operations in the top half of the kernel to register an unsatisfied asynchronous poll or select event request with the kernel. This registration enables later calls to the **selnotify** kernel service from resources in the bottom half of the kernel to correctly identify processes awaiting events on those resources.

The event requests may originate from calls to the **poll** or **select** subroutine, from processes, or from calls to the **fp\_poll** or **fp\_select** kernel service. A **select** file operation calls the **selreg** kernel service under the following circumstances:

- The poll or select request is asynchronous (the **POLLSYNC** flag is not set for the requested event's bit flags).
- The poll or select request determines (by calling the underlying resource's **ddselect** entry point) that the requested events have not yet occurred.

A registered event request takes the form of a select control block. The select control block is a structure containing the following:

- Requested event bit flags
- Returned event bit flags
- Primary resource identifier
- Unique resource identifier
- Pointer to a **proc** table entry
- File descriptor correlator
- Pointer to a notification routine that is non-null only for nested calls to the **poll** and **select** subroutines

The **selreg** kernel service allocates and initializes a select control block each time it is called.

When an event occurs on a resource that supports the **select** file operation, the resource calls the **selnotify** kernel service. The **selnotify** kernel service locates all select control blocks whose primary and unique identifiers match those of the resource, and whose requested event flags match the occurred events on the resource. Then, for each of the matching control blocks, the **selnotify** kernel service takes one of two courses of action, depending upon whether the control block's notification routine pointer is non-null (nested) or null (non-nested):

- In nested calls to the **select** or **poll** subroutines, the notification routine is called with the primary and unique resource identifiers, the returned event bit flags, and the process identifiers.
- In non-nested calls to the **select** or **poll** subroutine (the usual case), the SSEL bit of the process identified in the block is cleared, the returned event bit flags in the block are updated, and the process is awakened. A process awakened in this manner completes the **poll** or **select** call in which it was sleeping. The **poll** or **select** subroutine then collects the returned event bit flags in its processes' select control blocks for return to the user mode process, deallocates the control blocks, and returns tallies of the numbers of requested events that occurred to the user process.

## Execution Environment

The **selreg** kernel service can be called from the process environment only.

## Returns Values

Item	Description
0	Indicates successful completion.
EAGAIN	Indicates the <b>selreg</b> kernel service was unable to allocate a select control block.

### Related reference:

“**ddselect** Device Driver Entry Point” on page 633

“**fp\_select** Kernel Service” on page 167

### Related information:

select subroutine

Kernel Extension and Device Driver Management Kernel Services

## set\_pag or set\_pag64 Kernel Service Purpose

Sets a Process Authentication Group (PAG) value for the current process.

### Syntax

```
#include <sys/cred.h>
```

```
int set_pag ( type, pag )
int type;
int pag;
```

```
int set_pag64 ( type, pag )
int type;
uint64_t *pag;
```

### Parameters

Item	Description
<i>type</i>	PAG type to change
<i>pag</i>	PAG value

### Description

The **set\_pag** or **set\_pag64** kernel service copies the requested PAG for the current process. The caller must synchronize the **set\_pag** and **set\_pag64** kernel services with **validate\_pag** because **set\_pag** and **set\_pag64** do not lock process creation across the system. The value of *type* must be a defined PAG ID. The PAG ID for the Distributed Computing Environment (DCE) is 0.

### Execution Environment

The **set\_pag** and **set\_pag64** kernel services can be called from the process environment only.

### Return Values

A value of 0 is returned upon successful completion. Upon failure, a -1 is returned and **errno** is set to a value that explains the error.

### Error Codes

The **set\_pag** and **set\_pag64** kernel services fails if one or both of the following conditions are true:

Item	Description
EINVAL	Invalid PAG specification

### Related information:

Security Kernel Services

## setioctlr Subroutine Purpose

Sets a value to be returned by an **ioctl** routine.

### Syntax

```
void setioctlr ( ioctlr )
int ioctlr;
```

## Parameters

Item	Description
<i>ioctlrv</i>	Specifies an integer value to be returned by a successful completion of the <b>ioctl</b> subroutine.

## Description

The **setioclrv** subroutine sets the value of the `u_ioctlrv` field in the **uthread** structure of the running thread. The value in the `u_ioctlrv` field is returned by the **ioctl** or **fp\_ioctl** subroutine on a successful completion. If the **ioctl** subroutine fails, an `errno` value is returned instead.

## Return Values

The **setioclrv** subroutine returns no return values.

## Error Codes

The **setioclrv** subroutine returns no error codes.

## setjmpx Kernel Service

### Purpose

Allows saving the current execution state or context.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
int setjmpx ( jump_buffer)
label_t *jump_buffer;
```

### Parameter

Item	Description
<i>jump_buffer</i>	Specifies the address of the caller-supplied jump buffer that was specified on the call to the <b>setjmpx</b> service.

## Description

The **setjmpx** kernel service saves the current execution state, or context, so that a subsequent **longjmpx** call can cause an immediate return from the **setjmpx** service. The **setjmpx** service saves the context with the necessary state information including:

- The current interrupt priority.
- Whether the process currently owns the kernel mode lock.

Other state variables include the nonvolatile general purpose registers, the current program's table of contents and stack pointers, and the return address.

Calls to the **setjmpx** service can be nested. Each call to the **setjmpx** service causes the context at this point to be pushed to the top of the stack of saved contexts.

## Execution Environment

The **setjmpx** kernel service can be called from either the process or interrupt environment.

## Return Values

Item	Description
Nonzero value	Indicates that a <b>longjmp</b> call caused the <b>setjmp</b> service to return.
0	Indicates any other circumstances.

### Related reference:

“**clrjmp** Kernel Service” on page 43

### Related information:

Handling Signals While in a System Call

Exception Processing

## setpinit Kernel Service

### Purpose

Sets the parent of the current kernel process to the initialization process.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/device.h>
int setpinit()
```

### Description

The **setpinit** kernel service can be called by a kernel process to set its parent process to the **init** process. This is done to redirect the death of child signal for the termination of the kernel process. As a result, the **init** process is allowed to perform its default zombie process cleanup.

The **setpinit** service is used by a kernel process that can terminate, but does not want the user-mode process under which it was created to receive a death of child process notification.

### Execution Environment

The **setpinit** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
EINVAL	Indicates that the current process is not a kernel process.

### Related information:

Using Kernel Processes

Process and Exception Management Kernel Services

## setuerror Kernel Service

### Purpose

Allows kernel extensions to set the **ut\_error** field for the current thread.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
int setuerror ( errno)
int errno;
```

## Parameter

Item	Description
<i>errno</i>	Contains a value found in the <code>/usr/include/sys/errno.h</code> file that is to be copied to the current thread <code>ut_error</code> field.

## Description

The `setuerror` kernel service allows a kernel extension in a process environment to set the `ut_error` field in current thread's `uthread` structure. Kernel extensions providing system calls available to user-mode applications typically use this service. For system calls, the value of the `ut_error` field in the per thread `uthread` structure is copied to the `errno` global variable by the system call handler before returning to the caller.

## Execution Environment

The `setuerror` kernel service can be called from the process environment only.

## Return Codes

The `setuerror` kernel service returns the *errno* parameter.

### Related reference:

“getuerror Kernel Service” on page 191

### Related information:

Kernel Extension and Device Driver Management Kernel Services  
Understanding System Call Execution

## shutdown\_notify\_reg Kernel Service

### Purpose

Allows kernel extensions to register a shutdown notification.

### Syntax

```
#include <sys/reboot.h>

int shutdown_notify_reg(sn)
shutdown_notify_t *sn;

typedef struct _shutdown_notify {
    struct _shutdown_notify *next; /* Next in the link-list */
    int version; /* Version of structure */
    int oper; /* Bit map of the operation being performed */
    int status; /* The current status of this notify */
    uchar cb_retry; /* Internal use */
    uchar scope; /* Partition or system wide */
    uchar reason; /* User initiated or EPOW */
    uchar padding; /* padding */
    long (*func)(); /* Function kernel calls to notify ext. */
    void *uaddr;
} shutdown_notify_t;

/* Valid values
for shutdown_notify_t->oper */
#define SHUTDOWN_NOTIFY_PREPARE 0x1 /* Shutdown has started */
#define SHUTDOWN_NOTIFY_REBOOT 0x2 /*
```

```

Final notify that shutdown will be a reboot */
#define SHUTDOWN_NOTIFY_HALT 0x4
/* Final notify that shutdown will be a halt */
#define SHUTDOWN_NOTIFY_QUERY 0x8
/* Check to see if finished shutdown */

/* Valid values for
shutdown_notify_t->status and
for SHUTDOWN_NOTIFY_QUERY return code */
#define SHUTDOWN_STATUS_PREPARE 0x1 /* Preparing for shutdown */
#define SHUTDOWN_STATUS_COMMENCE 0x2 /* Commencing shutdown */
#define SHUTDOWN_STATUS_FINISH 0x4 /* Finished shutdown */

#define SHUTDOWN_NOTIFY_VERSION 1 /* Increment by 1
    * every time add more
    * variables to
    * shutdown_notify_t
    */
/* Valid values for shutdown_notify_t->scope */
#define SHUTDOWN_SCOPE_PARTITION 1
#define SHUTDOWN_SCOPE_SYSTEM 2

/* Valid values for shutdown_notify_t->reason */
#define SHUTDOWN_REASON_USER 1
#define SHUTDOWN_REASON_EPOW 2

/* Valid handler return codes
during the SHUTDOWN_NOTIFY_PREPARE phase */
#define SHUTDOWN_RC_SUCCESS 0
#define SHUTDOWN_RC_DELAY 1

#define SHUTDOWN_NOTIFY_VERSION 2

```

## Description

The shutdown notify subsystem has been extended to provide additional information during a shutdown operation. During the **SHUTDOWN\_NOTIFY\_PREPARE** phase, the kernel provides information on the scope and reason for the shutdown action. Additionally, when a handler is called, before its completion, it can now delay the shutdown operation in order to finalize any outstanding jobs. The kernel again then calls out to the handler after some small amount of time. This process continues until all handlers return **SHUTDOWN\_RC\_SUCCESS**. This functionality is only present for **shutdown\_notify\_t** version 2 and preceding handlers. For version 1 handlers, the new fields are not present and the return code from the handler is ignored.

## Parameters

Item	Description
<i>cb_retry</i>	Internal use.
<i>func</i>	Pointer to the function called to notify registered extension.
<i>next</i>	Pointer to next <b>shutdown_notify_t</b> structure in list.
<i>oper</i>	Bit map of operation(s) being performed.
<i>padding</i>	Padding.
<i>reason</i>	User initiated or EPOW event.
<i>scope</i>	Shutdown at the partition or system level.
<i>sn</i>	Pointer to a structure that the calling extension fills out when it registers.
<i>status</i>	Current status of notify.
<i>uaddr</i>	Place for extension to store an address to help it identify the object to which this structure refers.
<i>version</i>	Version of structure. Set to 1.
<b>SHUTDOWN_NOTIFY_HALT</b>	A halt is occurring.
<b>SHUTDOWN_NOTIFY_PREPARE</b>	Shutdown has started.
<b>SHUTDOWN_NOTIFY_QUERY</b>	Check to see if finished shutdown.

Item	Description
<i>SHUTDOWN_NOTIFY_REBOOT</i>	A reboot is occurring.
<i>SHUTDOWN_NOTIFY_VERSION</i>	Version number of structure.
<i>SHUTDOWN_RC_DELAY</i>	Return from registered handler to indicate its processing is not complete and wants to delay the shutdown operation.
<i>SHUTDOWN_RC_SUCCESS</i>	Return from registered handler to indicate all processing is complete and the shutdown operation can proceed.
<i>SHUTDOWN_REASON_EPOW</i>	EPOW event.
<i>SHUTDOWN_REASON_USER</i>	User initiated shutdown.
<i>SHUTDOWN_SCOPE_PARTITION</i>	Shutdown at the partition level.
<i>SHUTDOWN_SCOPE_SYSTEM</i>	Shutdown at the system level.
<i>SHUTDOWN_STATUS_COMMENCE</i>	Wrap up shutdown.
<i>SHUTDOWN_STATUS_FINISH</i>	Shutdown has completed.
<i>SHUTDOWN_STATUS_PREPARE</i>	Preparing for shutdown.

## Execution Environment

Process environment only.

## Return Values

Item	Description
0	Success.
EPERM	Attempted to register after prepare notification has started.
EINVAL	Invalid argument passed.

### Related reference:

“shutdown\_notify\_unreg Kernel Service”

## shutdown\_notify\_unreg Kernel Service

### Purpose

Unregisters an extension from getting notified in the event of a shutdown.

### Syntax

```
#include <sys/reboot.h>
```

```
int shutdown_notify_unreg(sn)
shutdown_notify_t *sn;
```

### Description

The **shutdown\_notify\_unreg** kernel service unregisters an extension from getting notified in the event of a shutdown. The extension passes in the **shutdown\_notify\_t** instance it wants to unregister. This function will fail if it is called after the **SHUTDOWN\_NOTIFY\_HALT** and **SHUTDOWN\_NOTIFY\_REBOOT** notification process has started.

### Parameters



Item	Description
<i>sn</i>	Pointer to a structure that the calling extension wants to unregister.

## Execution Environment

Process environment only.

## Return Values

Item	Description
0	Success
EPERM	Attempted to unregister after final notification has started.
EINVAL	Invalid argument passed.

### Related reference:

“shutdown\_notify\_reg Kernel Service” on page 470

## sig\_chk Kernel Service

### Purpose

Provides a kernel process the ability to poll for receipt of signals.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/signal.h>
int sig_chk ()
```

### Description

**Attention:** A system crash will occur if the **sig\_chk** service is not called by a kernel process.

The **sig\_chk** kernel service can be called by a kernel thread in kernel mode to determine if any unmasked signals have been received. Signals do not preempt threads because serialization of critical data areas would be lost. Instead, threads must poll for signals, either periodically or after a long sleep has been interrupted by a signal.

The **sig\_chk** service checks for any pending signal that has a specified *signal catch* or *default* action. If one is found, the service returns the signal number as its return value. It also removes the signal from the pending signal mask. If no signal is found, this service returns a value of 0. The **sig\_chk** service does not return signals that are blocked or ignored. It is the responsibility of the kernel process to handle the signal appropriately.

For kernel-only threads, the **sig\_chk** kernel service clears the returned signal from the list of pending signals. For other kernel threads, the signal is not cleared, but left pending. It will be delivered to the kernel thread as soon as it returns to the user mode.

Understanding Kernel Threads in *Kernel Extensions and Device Support Programming Concepts* provides more information about kernel-only thread signal handling.

## Execution Environment

The **sig\_chk** kernel service can be called from the process environment only.

## Return Values

Upon completion, the **sig\_chk** service returns a value of 0 if no pending unmasked signal is found. Otherwise, it returns a nonzero signal value indicating the number of the highest priority signal that is pending. Signal values are defined in the `/usr/include/sys/signal.h` file.

### Related information:

Introduction to Kernel Processes

Process and Exception Management Kernel Services

## **simple\_lock** or **simple\_lock\_try** Kernel Service Purpose

Locks a simple lock.

### Syntax

```
#include <sys/lock_def.h>
```

```
void simple_lock ( lock_addr)  
simple_lock_t lock_addr;
```

```
boolean_t simple_lock_try ( lock_addr)  
simple_lock_t lock_addr;
```

### Parameter

Item	Description
<i>lock_addr</i>	Specifies the address of the lock word to lock.

### Description

The **simple\_lock** kernel service locks the specified lock; it blocks if the lock is busy. The lock must have been previously initialized with the **simple\_lock\_init** kernel service. The **simple\_lock** kernel service has no return values.

The **simple\_lock\_try** kernel service tries to lock the specified lock; it returns immediately without blocking if the lock is busy. If the lock is free, the **simple\_lock\_try** kernel service locks it. The lock must have been previously initialized with the **simple\_lock\_init** kernel service.

**Note:** When using simple locks to protect thread-interrupt critical sections, it is recommended that you use the **disable\_lock** kernel service instead of calling the **simple\_lock** kernel service directly.

### Execution Environment

The **simple\_lock** and **simple\_lock\_try** kernel services can be called from the process environment only.

### Return Values

The **simple\_lock\_try** kernel service has the following return values:

Item	Description
TRUE	Indicates that the simple lock has been successfully acquired.
FALSE	Indicates that the simple lock is busy, and has not been acquired.

**Related reference:**

“disable\_lock Kernel Service” on page 77

“simple\_unlock Kernel Service” on page 476

**Related information:**

Understanding Locking

Locking Kernel Services

## simple\_lock\_init Kernel Service Purpose

Initializes a simple lock.

### Syntax

```
#include <sys/lock_def.h>
```

```
void simple_lock_init ( lock_addr)
simple_lock_t lock_addr;
```

### Parameter

Item	Description
<i>lock_addr</i>	Specifies the address of the lock word.

### Description

The **simple\_lock\_init** kernel service initializes a simple lock. This kernel service must be called before the simple lock is used. The simple lock must previously have been allocated with the **lock\_alloc** kernel service.

### Execution Environment

The **simple\_lock\_init** kernel service can be called from the process environment only.

The **simple\_lock\_init** kernel service may be called either the process or interrupt environments.

### Return Values

The **simple\_lock\_init** kernel service has no return values.

**Related reference:**

“lock\_alloc Kernel Service” on page 338

“simple\_lock or simple\_lock\_try Kernel Service” on page 474

“simple\_unlock Kernel Service” on page 476

**Related information:**

Understanding Locking

Locking Kernel Services

## simple\_unlock Kernel Service

### Purpose

Unlocks a simple lock.

### Syntax

```
#include <sys/lock_def.h>
```

```
void simple_unlock ( lock_addr)  
simple_lock_t lock_addr;
```

### Parameter

Item	Description
<i>lock_addr</i>	Specifies the address of the lock word to unlock.

### Description

The **simple\_unlock** kernel service unlocks the specified simple lock. The lock must be held by the thread which calls the **simple\_unlock** kernel service. Once the simple lock is unlocked, the highest priority thread (if any) which is waiting for it is made runnable, and may compete for the lock again. If at least one kernel thread was waiting for the lock, the priority of the calling kernel thread is recomputed.

**Note:** When using simple locks to protect thread-interrupt critical sections, it is recommended that you use the **unlock\_enable** kernel service instead of calling the **simple\_unlock** kernel service directly.

### Execution Environment

The **simple\_unlock** kernel service can be called from the process environment only.

### Return Values

The **simple\_unlock** kernel service has no return values.

#### Related reference:

“simple\_lock\_init Kernel Service” on page 475

“simple\_lock or simple\_lock\_try Kernel Service” on page 474

“unlock\_enable Kernel Service” on page 518

#### Related information:

Understanding Locking

## sleep Kernel Service

### Purpose

Forces the calling kernel thread to wait on a specified channel.

### Syntax

```
#include <sys/types.h>  
#include <sys/errno.h>  
#include <sys/pri.h>  
#include <sys/proc.h>
```

```
int sleep ( chan, priflags)  
void *chan;  
int priflags;
```

## Parameters

Item	Description
<i>chan</i>	Specifies the channel number. For the <b>sleep</b> service, this parameter identifies the channel to wait for (sleep on).
<i>priflags</i>	Specifies two conditions: <ul style="list-style-type: none"><li>• The priority at which the kernel thread is to run when it is reactivated.</li><li>• Flags indicating how a signal is to be handled by the <b>sleep</b> kernel service.</li></ul>

The valid flags and priority values are defined in the `/usr/include/sys/pri.h` file.

## Description

The **sleep** kernel service is provided for compatibility only and should not be invoked by new code. The **e\_sleep\_thread** or **et\_wait** kernel service should be used when writing new code.

The **sleep** service puts the calling kernel thread to sleep, causing it to wait for a wakeup to be issued for the channel specified by the *chan* parameter. When the process is woken up again, it runs with the priority specified in the *priflags* parameter. The new priority is effective until the process returns to user mode.

All processes that are waiting on the channel are restarted at once, causing a race condition to occur between the activated threads. Thus, after returning from the **sleep** service, each thread should check whether it needs to sleep again.

The channel specified by the *chan* parameter is simply an address that by convention identifies some event to wait for. When the kernel or kernel extension detects such an event, the **wakeup** service is called with the corresponding value in the *chan* parameter to start up all the threads waiting on that channel. The channel identifier must be unique systemwide. The address of an external kernel variable (which can be defined in a device driver) is generally used for this value.

If the **SWAKEONSIG** flag is not set in the *priflags* parameter, signals do not terminate the sleep. If the **SWAKEONSIG** flag is set and the **PCATCH** flag is not set, the kernel calls the **longjmpx** kernel service to resume the context saved by the last **setjmpx** call if a signal interrupts the sleep. Therefore, any system call (such as those calling device driver **ddopen**, **ddread**, and **ddwrite** routines) or kernel process that does an interruptible sleep without the **PCATCH** flag set must have set up a context using the **setjmpx** kernel service. This allows the sleep to resume in case a signal is sent to the sleeping process.

**Attention:** The caller of the **sleep** service must own the kernel-mode lock specified by the *kernel\_lock* parameter. The **sleep** service does not provide a compatible level of serialization if the kernel lock is not owned by the caller of the **sleep** service.

## Execution Environment

The **sleep** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates successful completion.
1	Indicates that a signal has interrupted a sleep with both the <b>PCATCH</b> and <b>SWAKEONSIG</b> flags set in the <i>priflags</i> parameter.

#### Related information:

Locking Strategy in Kernel Mode

Process and Exception Management Kernel Services

### | **sleepx Kernel Service**

#### | **Purpose**

| Wait for an event.

#### | **Syntax**

| #include <sys/sleep.h>

| **int sleepx** (*tchan\_t chan int pri flags\_t flags*)

#### | **Parameters**

##### | **chan**

| Specifies the channel number. For the **sleep** service, this parameter identifies the channel to wait for (sleep on).

##### | **pri**

| Specifies the wakeup priority

##### | **flags**

| Signal control flags

#### | **Description**

| Wait for an event to occur. This procedure can only be called by a thread. Callers of this service must be prepared for a premature return and check that the reason for waiting has gone away.

| The **pri** parameter will be the priority of the thread when it becomes runnable again (if that priority is more favorable). The process will keep that priority until it is dispatched. The range of the wakeup priority is  $0 \leq pri \leq PRI\_LOW$ . If the **pri** parameter is outside of that range, it is forced to the lower or upper boundary.

#### | **Execution Environment**

| The **sleepx** kernel service can be called from the process environment only.

#### | **Return Values**

| 0 Indicates that the event occurred.

| 1 Indicates that the event signalled out.

### | **subyte Kernel Service**

#### **Purpose**

Stores a byte of data in user memory.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
int subyte ( uaddr, c )
uchar *uaddr;
uchar c;
```

## Parameters

Item	Description
<i>uaddr</i>	Specifies the address of user data.
<i>c</i>	Specifies the character to store.

## Description

The **subyte** kernel service stores a byte of data at the specified address in user memory. It is provided so that system calls and device heads can safely access user data. The **subyte** service ensures that the user has the appropriate authority to:

- Access the data.
- Protect the operating system from paging I/O errors on user data.

The **subyte** service should only be called while executing in kernel mode in the user process.

## Execution Environment

The **subyte** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates successful completion.
-1	Indicates a <i>uaddr</i> parameter that is not valid for one of the following reasons: <ul style="list-style-type: none"><li>• The user does not have sufficient authority to access the data.</li><li>• The address is not valid.</li><li>• An I/O error occurs when the user data is referenced.</li></ul>

### Related reference:

“fubyte Kernel Service” on page 179

### Related information:

Accessing User-Mode Data While in Kernel Mode  
Memory Kernel Services

## suser Kernel Service Purpose

Determines the privilege state of a process.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
int suser ( ep )
char *ep;
```

## Parameter

Item	Description
<i>ep</i>	Points to a character variable where the <b>EPERM</b> value is stored on failure.

## Description

The **suser** kernel service checks whether a process has any effective privilege (that is, whether the process's uid field equals 0).

## Execution Environment

The **suser** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates failure. The character pointed to by the <i>ep</i> parameter is set to the value of <b>EPERM</b> . This indicates that the calling process does not have any effective privilege.
Nonzero value	Indicates success (the process has the specified privilege).

### Related information:

Security Kernel Services

## suword Kernel Service

### Purpose

Stores a word of data in user memory.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
int suword ( uaddr, w)
int *uaddr;
int w;
```

### Parameters

Item	Description
<i>uaddr</i>	Specifies the address of user data.
<i>w</i>	Specifies the word to store.

## Description

The **suword** kernel service stores a word of data at the specified address in user memory. It is provided so that system calls and device heads can safely access user data. The **suword** service ensures that the user had the appropriate authority to:

- Access the data.
- Protect the operating system from paging I/O errors on user data.

The **suword** service should only be called while executing in kernel mode in the user process.



## Execution Environment

The `suword` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates successful completion.
-1	Indicates a <code>uaddr</code> parameter that is not valid for one of these reasons: <ul style="list-style-type: none"><li>• The user does not have sufficient authority to access the data.</li><li>• The address is not valid.</li><li>• An I/O error occurs when the user data is referenced.</li></ul>

### Related reference:

“`fuword` Kernel Service” on page 180

### Related information:

Memory Kernel Services

Accessing User-Mode Data While in Kernel Mode

## t

The following kernel services begin with the with the letter t.

## TE\_verify\_reg Kernel Service

### Purpose

Registers a callout handler for Trusted Execution (TE) file verification during the `exec()` functions, kernel extension loads, and library load operations.

### Syntax

```
#include <sys/file.h>
typedef int (*TE_verify)(char *, int, struct file *);

int TE_verify_reg(TE_verify verify_fn, uint_64 options)
```

### Parameters

#### **verify\_fn**

Specifies the callout function to be called for the verification checks with the `exec()` functions for the Trusted Execution of the AIX kernel level, loading of kernel extensions, and library loading events instead of the default AIX Trusted Execution method.

For more information about the function definition of this callout handler, see the `alt_verify_fn` section.

#### **options**

Specifies a bit mask of registration options. The **options** parameter is not defined currently. The caller must set the **options** parameter to 0.

### Description

The `TE_verify_reg` kernel service registers a callout handler for the AIX Trusted Execution framework.

After a callout handler is registered, the handler is invoked for the `exec()` functions, loading kernel extensions, and library load-time checks for Trusted Execution in the AIX kernel. The default AIX Trusted Execution logic is not invoked and any AIX-configured policies for Trusted Execution not applied. The registered alternative handler becomes the active Trusted Execution engine for AIX to provide security policy as implemented in the handler and its associated management components.

After a callout handler is registered with the `TE_verify_reg` kernel service, subsequent invocation of the `TE_verify_reg` service returns with an error code of `EEXIST`.

You must have root authority to call the `TE_verify_reg` kernel service.

## Return values

On successful completion, the `TE_verify_reg` service kernel service returns a value of 0.

The following error codes are returned on failure:

### **EEXIST**

The callout handler is already registered.

### **EPERM**

The caller does not have permission to invoke this function.

### **EINVAL**

The callout handler or the **options** parameters are invalid.

## Execution environment

The `TE_verify_reg` kernel service can be called from the process environment only.

The registered alternative Trusted Execution handler must conform to the behaviors that are described in the following section.

## alt\_verify\_fn callout function

### Purpose

Verifies the integrity of a file.

### Syntax

```
#include <sys/file.h>

#define VERIFY_EXECUTABLES 2
#define VERIFY_SHLIBS 3
#define VERIFY_SCRIPTS 4
#define VERIFY_KERNEXTS 5

int alt_TE_verify (char *path_name, int type, struct file *path_fp)
```

### Description

The `alt_TE_verify` callout function is started from the loader and the program execution path to verify the integrity of a file that is specified under the `path_name` parameter. The `path_fp` parameter is a file pointer to the file object that is associated with the `path_name` parameter.

The `type` parameter can be one of the following values:

#### **VERIFY\_EXECUTABLES**

This value is specified when the `alt_TE_verify` function is started from the kernel `exec()` function to verify executable programs.

#### **VERIFY\_SCRIPTS**

This value is specified when the `alt_TE_verify` function is started from the `exec()` function and the `path_name` value is a shell file.

## VERIFY\_KERNEXTS

This value is specified when the `alt_TE_verify` function is started for loading a kernel extension.

## VERIFY\_SHLIBS

This value is specified when the `alt_TE_verify` function is started for loading a shared library.

## Input parameters

### `path_name`

Specifies the path to the file that must be verified.

### `type`

Indicates the type of verification that must be performed.

### `path_fp`

Indicates the file pointer to the `path_name` file.

## Return values

**0** Indicates that the verification completed successfully.

### **Nonzero**

Indicates that the verification failed.

The nonzero return value blocks loading of the file. An error number is set by the AIX kernel functions that start the `alt_verify_fn` callout function.

## TE\_verify\_unreg Kernel Service

### Purpose

Unregisters a previously registered callout handler for trusted execution.

### Syntax

```
#include <sys/file.h>
typedef int (*TE_verify)(char *, int, struct file *);

int TE_verify_unreg(TE_verify verify_fn, uint_64 options)
```

### Parameters

#### `verify_fn`

Specifies the callout function that must be used when you register the handler by using the `TE_verify_reg()` kernel service.

#### `options`

Specifies a bit mask of registration options. The **options** parameter is not defined currently. The caller must set the **options** parameter to 0.

### Description

The `TE_verify_unreg` kernel service unregisters a callout handler for the AIX Trusted Execution (TE) framework. The `verify_fn` parameter must match with the currently registered TE callout handler. Otherwise, the `TE_verify_unreg` kernel service returns an error code of `EPERM`.

After a callout handler is unregistered, the default AIX trusted execution logic is applied based on the configured AIX trusted execution policies.

The caller of the `TE_verify_unreg` kernel service must have root authority.

## Return values

On successful completion, the `TE_verify_unreg` kernel service returns a value of 0.

The following error codes are returned on failure:

### EPERM

The caller does not have permission to start this function. Or, the registered callout handler is not same as the *verify\_fn* parameter.

### EINVAL

No callout handler is registered or the *options* parameters are invalid.

## Execution environment

The `TE_verify_unreg` kernel service can be called only from the process environment.

## talloc Kernel Service

### Purpose

Allocates a timer request block before starting a timer request.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/timer.h>
struct trb *talloc()
```

### Description

The `talloc` kernel service allocates a timer request block. The user must call it before starting a timer request with the `tstart` kernel service. If successful, the `talloc` service returns a pointer to a pinned timer request block.

### Execution Environment

The `talloc` kernel service can be called from the process environment only.

### Return Values

The `talloc` service returns a pointer to a timer request block upon successful allocation of a `trb` structure. Upon failure, a null value is returned.

#### Related reference:

“tfree Kernel Service”

#### Related information:

Timer and Time-of-Day Kernel Services

Using Fine Granularity Timer Services and Structures

## tfree Kernel Service

### Purpose

Deallocates a timer request block.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/timer.h>
```

```
void tfree ( t )
struct trb *t;
```

## Parameter

Item	Description
<i>t</i>	Points to the timer request structure to be freed.

## Description

The **tfree** kernel service deallocates a timer request block that was previously allocated with a call to the **talloc** kernel service. The caller of the **tfree** service must first cancel any pending timer request associated with the timer request block being freed before attempting to free the request block. Canceling the timer request block can be done using the **tstop** kernel service.

## Execution Environment

The **tfree** kernel service can be called from either the process or interrupt environment.

**Note:** Do not use the **tfree** kernel service to free the timer request block that is passed to the timer completion handler.

## Return Values

The **tfree** service has no return values.

### Related reference:

“talloc Kernel Service” on page 484

### Related information:

Timer and Time-of-Day Kernel Services

Using Fine Granularity Timer Services and Structures

## thread\_create Kernel Service

### Purpose

Creates a new kernel thread in the calling process.

### Syntax

```
#include <sys/thread.h>
tid_t thread_create ()
```

### Description

The **thread\_create** kernel service creates a new kernel-only thread in the calling kernel process. The thread's ID is returned; it is unique system wide.

The new thread does not begin running immediately; its state is set to **TSIDL**. The execution will start after a call to the **kthread\_start** kernel service. If the process is exited prior to the thread being made

runnable, the thread's resources are released immediately. The thread's signal mask is inherited from the calling thread; the set of pending signals is cleared. Signals sent to the thread are marked pending while the thread is in the **TSIDL** state.

If the calling thread is bound to a specific processor, the new thread will also be bound to the processor.

## Execution Environment

The **thread\_create** kernel service can be called from the process environment only. This service cannot be called directly from a kernel extension.

## Return Values

Upon successful completion, the new thread's ID is returned. Otherwise, -1 is returned, and the error code can be checked by calling the **getuerror** kernel service.

## Error Codes

Item	Description
EAGAIN	The total number of kernel threads executing system wide or the maximum number of kernel threads per process would be exceeded.
ENOMEM	There is not sufficient memory to create the kernel thread.
ENOTSUP	The <b>thread_create</b> service was called directly from a kernel extension.

### Related reference:

"kthread\_start Kernel Service" on page 311

### Related information:

Process and Exception Management Kernel Services

## thread\_self Kernel Service

### Purpose

Returns the caller's kernel thread ID.

### Syntax

```
#include <sys/thread.h>
tid_t thread_self ()
```

### Description

The **thread\_self** kernel service returns the thread process ID of the calling process.

The **thread\_self** service can also be used to check the environment that the routine is being executed in. If the caller is executing in the interrupt environment, the **thread\_self** service returns a process ID of -1. If a routine is executing in a process environment, the **thread\_self** service obtains the thread process ID.

## Execution Environment

The **thread\_self** kernel service can be called from either the process or interrupt environment.

## Return Values

Item	Description
-1	Indicates that the <code>thread_self</code> service was called from an interrupt environment.

The `thread_self` service returns the thread process ID of the current process if called from a process environment.

#### Related information:

Process and Exception Management Kernel Services  
Understanding Execution Environments

## thread\_setsched Kernel Service

### Purpose

Sets kernel thread scheduling parameters.

### Syntax

```
#include <sys/thread.h>
#include <sys/sched.h>
```

```
int thread_setsched ( tid, priority, policy)
tid_t tid;
int priority;
int policy;
```

### Parameters

Item	Description
<i>tid</i>	Specifies the kernel thread.
<i>priority</i>	Specifies the priority. It must be in the range from 0 to <code>PRI_LOW</code> ; 0 is the most favored priority.
<i>policy</i>	Specifies the scheduling policy. It must have one of the following values:
<code>SCHED_FIFO</code>	Denotes fixed priority first-in first-out scheduling.
<code>SCHED_FIFO2</code>	Allows a thread that sleeps for a relatively short amount of time to be requeued to the head, rather than the tail, of its priority run queue.
<code>SCHED_FIFO3</code>	Causes threads to be enqueued to the head of their run queues.
<code>SCHED_RR</code>	Denotes fixed priority round-robin scheduling.
<code>SCHED_OTHER</code>	Denotes the default scheduling policy.

### Description

The `thread_setsched` subroutine sets the scheduling parameters for a kernel thread. This includes both the priority and the scheduling policy, which are specified in the *priority* and *policy* parameters. The calling and the target thread must be in the same process.

When setting the scheduling policy to `SCHED_OTHER`, the system chooses the priority; the *priority* parameter is ignored. The only way to influence the priority of a thread using the default scheduling policy is to change the process nice value.

The calling thread must belong to a process with root authority to change the scheduling policy of a thread to either `SCHED_FIFO`, `SCHED_FIFO2`, `SCHED_FIFO3`, or `SCHED_RR`.

## Execution Environment

The `thread_setsched` kernel service can be called from the process environment only.

## Return Values

Upon successful completion, 0 is returned. Otherwise, -1 is returned, and the error code can be checked by calling the `getuerror` kernel service.

## Error Codes

Item	Description
EINVAL	The <i>priority</i> or <i>policy</i> parameters are not valid.
EPERM	The calling kernel thread does not have sufficient privilege to perform the operation.
ESRCH	The kernel thread <i>tid</i> does not exist.

### Related reference:

“thread\_create Kernel Service” on page 485

### Related information:

Process and Exception Management Kernel Services

## thread\_set\_smt\_priority or thread\_read\_smt\_priority System Call Purpose

Sets or reads the current simultaneous multithreading (SMT) thread priority for a user-thread.

## Syntax

```
#include <sys/errno.h>
#include <sys/thread.h>
#include <sys/processor.h>
```

```
int thread_set_smt_priority ( Priority )
smt_thread_priority_t Priority;
```

```
#include <sys/errno.h>
#include <sys/thread.h>
#include <sys/processor.h>
```

```
smt_thread_priority_t thread_read_smt_priority ( )
```

## Description

The SMT thread priority that is associated with a logical CPU, SMT hardware thread, controls the relative priority of the logical CPU in relation to the other logical CPUs on the same processor core. The relative priority between the SMT hardware threads on a processor core determines how decode cycles are granted to each SMT hardware thread. The SMT thread priority can be used to cause a particular application thread to be favored over other application threads that are running on the other SMT hardware threads in the same processor core. It is done by increasing the SMT thread priority of the logical CPU the application is running on, or by lowering the SMT thread priority of the application threads that are running on the other logical CPUs associated with the same processor core.

The `thread_set_smt_priority` and `thread_read_smt_priority` system calls provide a way to register and read back the current SMT thread priority on a per process-thread basis.

### Note:

These interfaces are not supported on some processor architectures.



If the process-thread is dispatched to a logical CPU that is running in non-SMT mode, the SMT thread priority level has no effect.

Callers of the **thread\_set\_smt\_priority** system call with normal user-level privileges can set their SMT thread priority level to one of the following levels:

- LOW
- MEDIUM LOW
- NORMAL

Callers that have RBAC PV\_PROC\_VARS privilege can set their priority level to one of the following levels:

- VERY LOW
- LOW
- MEDIUM LOW
- NORMAL
- MEDIUM HIGH
- HIGH

The default thread priority level is NORMAL.

**Note:** The only supported means for altering the SMT thread priority level is by using the **thread\_set\_smt\_priority** system call. If an alternative means of setting the SMT priority is used, the kernel does not know the process-thread's current SMT priority level, and overwrites the required SMT priority level without restoring it.

The **thread\_read\_smt\_priority** system call returns the current SMT priority level that is registered by the process thread. If the process thread did not register a required SMT priority level, then the default priority level of NORMAL is returned.

## Parameters

Item	Description
Priority	Used to specify one of the following parameters: <ul style="list-style-type: none"><li>• T_VERYLOW_SMT_PRI</li><li>• T_LOW_SMT_PRI</li><li>• T_MEDIUMLOW_SMT_PRI</li><li>• T_NORMAL_SMT_PRI</li><li>• T_MEDIUMHIGH_SMT_PRI</li><li>• T_HIGH_SMT_PRI</li></ul>

## Execution Environment

The **thread\_read\_smt\_priority** and **thread\_set\_smt\_priority** system calls can be called from the process environment only.

## Return Values

On successful completion, the **thread\_set\_smt\_priority** system call returns 0. Otherwise, **-1** is returned and the **errno** global variable is set to indicate the error.

On successful completion, the **thread\_read\_smt\_priority** system call returns the current required SMT priority. Otherwise, **-1** is returned and the **errno** global variable is set to indicate the error.

## Error Codes

Item	Description
EPERM	The process attempted to set the SMT thread priority level to a value other than T_LOW_SMT_PRI, T_MEDIUMLOW_SMT_PRI, or T_NORMAL_SMT_PRI and does not have the necessary privileges.
EINVAL	The required priority value that is specified is invalid.
ENOSYS	SMT thread priority level manipulation is not supported on this system.

## thread\_terminate Kernel Service

### Purpose

Terminates the calling kernel thread.

### Syntax

```
#include <sys/thread.h>
void thread_terminate ()
```

### Description

The **thread\_terminate** kernel service terminates the calling kernel thread and cleans up its structure and its kernel stack. If it is the last thread in the process, the process will exit.

The **thread\_terminate** kernel service is automatically called when a thread returns from its entry point routine (defined in the call to the **kthread\_start** kernel service).

### Execution Environment

The **thread\_terminate** kernel service can be called from the process environment only.

### Return Values

The **thread\_terminate** kernel service never returns.

#### Related reference:

“kthread\_start Kernel Service” on page 311

#### Related information:

Process and Exception Management Kernel Services

## timeout Kernel Service

**Attention:** This service must not be used because it is not multi-processor safe. The base kernel timer and watchdog services must be used instead.

### Purpose

Schedules a function to be called after a specified interval.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
void timeout ( func, arg, ticks)
void (*func)();
caddr_t *arg;
int ticks;
```

## Parameters

Item	Description
<i>func</i>	Indicates the function to be called.
<i>arg</i>	Indicates the parameter to supply to the function specified by the <i>func</i> parameter.
<i>ticks</i>	Specifies the number of timer ticks that must occur before the function specified by the <i>func</i> parameter is called. Many timer ticks can occur per second. The HZ label that is found in the <code>/usr/include/sys/m_param.h</code> file can be used to determine the number of ticks per second.

## Description

The **timeout** service is not part of the kernel. However, it is a compatibility service that is provided in the **libsys.a** library. To use the **timeout** service, a kernel extension must be bound with the **libsys.a** library. The **timeout** service, like the associated kernel services **untimeout** and **timeoutcf**, can be bound and used only in the pinned part of a kernel extension or the bottom half of a device driver because these services use interrupt disable for serialization.

The **timeout** service schedules the function pointed to by the *func* parameter to be called with the *arg* parameter after the number of timer ticks that are specified by the *ticks* parameter. Use the **timeoutcf** routine to allocate enough callout elements for the maximum number of simultaneous active time outs that you expect.

**Note:** The **timeoutcf** routine must be called before the **timeout** service is called.

Calling the **timeout** service without allocating enough callout table entries can result in a kernel panic because of a lack of pinned callout table elements. The value of a timer tick depends on the hardware's capability. You can use the **restimer** subroutine to determine the minimum granularity.

Multiple pending **timeout** requests with the same *func* and *arg* parameters are not allowed.

### The *func* Parameter

The function that is specified by the *func* parameter must be declared as follows:

```
void func (arg)
void *arg;
```

## Execution Environment

The **timeout** routine can be called from either the process or interrupt environment.

The function that is specified by the *func* parameter is called in the interrupt environment. Therefore, it must follow the conventions for interrupt handlers.

## Return Values

The **timeout** service has no return values.

### Related reference:

“untimeout Kernel Service” on page 523

“timeoutcf Subroutine for Kernel Services” on page 492

### Related information:

restimer subroutine

Timer and Time-of-Day Kernel Services

## timeoutcf Subroutine for Kernel Services

**Attention:** This service must not be used because it is not multi-processor safe. The base kernel timer and watchdog services must be used instead.

### Purpose

Allocates or deallocates callout table entries for use with the **timeout** kernel service.

### Library

libsys.a (Kernel extension runtime routines)

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
int timeoutcf ( cocnt )
int cocnt ;
```

### Parameter

Item	Description
<i>cocnt</i>	Specifies the callout count. This value indicates the number of callout elements by which to increase or decrease the current allocation. If this number is positive, the number of callout entries for use with the <b>timeout</b> service is increased. If this number is negative, the number of elements is decreased by the amount specified.

### Description

The **timeoutcf** subroutine is not part of the kernel. It is a compatibility service that is provided in the **libsys.a** library. To use the **timeoutcf** subroutine, a kernel extension must be bound with the **libsys.a** library. The **timeoutcf** subroutine, like the associated kernel **libsys** services **untimeout** and **timeout**, can be bound and used only in the pinned part of a kernel extension or the bottom half of a device driver because these services use interrupt disable for serialization.

The **timeoutcf** subroutine registers an increase or decrease in the number of callout table entries available for the **timeout** subroutine to use. Before a subroutine can use the **timeout** kernel service, the **timeoutcf** subroutine must increase the number of callout table entries available to the **timeout** kernel service. It increases this number by the maximum number of outstanding time outs that the routine can have pending at one time.

The **timeoutcf** subroutine must be used to decrease the number of callout table entries by the amount it was increased under the following conditions:

- The routine that uses the **timeout** subroutine finished using it.
- The calling routine has no more outstanding timeout requests that are pending.

Typically the **timeoutcf** subroutine is called in a device driver's **open** and **close** routine. It is called to allocate and deallocate sufficient elements for the maximum expected use of the **timeout** kernel service for that instance of the open device.

**Attention:** A kernel panic results either of these two circumstances:

- A request to decrease the callout table allocation is made that is greater than the number of unused callout table entries.
- The **timeoutcf** subroutine is called in an interrupt environment.

## Execution Environment

The `timeoutcf` subroutine can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful allocation or deallocation of the requested callout table entries.
-1	Indicates an unsuccessful operation.

### Related reference:

“timeout Kernel Service” on page 490

### Related information:

Timer and Time-of-Day Kernel Services

## trc\_ishookon Exported Kernel Service

### Purpose

Checks if a given trace hook word is being traced by system trace.

### Syntax

```
#include <sys/trcmacros.h>
```

```
int trc_ishookon (int chan, long hkwd);
```

### Description

The `trc_ishookon` kernel service informs the user if tracing is on and the specified hook word is being traced.

### Parameters

Item	Description
<i>chan</i>	The channel to query with the range from 0 to 7.
<i>hkwd</i>	The hook word to be traced by system trace.

### Return Values

Item	Description
1	The hook word is being traced.
0	Hook word is not being traced or system trace is off.

### Related information:

trace subroutine

## trcgenk Kernel Service

### Purpose

Records a trace event for a generic trace channel.

### Syntax

```
#include <sys/types.h>  
#include <sys/errno.h>  
#include <sys/trchkid.h>
```

```
void trcgenk (chan, hk_word, data_word, len, buf)
unsigned int chan, hk_word, data_word, len;
char * buf;
```

## Parameters

Item	Description				
<i>chan</i>	Specifies the channel number for the trace session. This number is obtained from the <b>trcstart</b> subroutine.				
<i>hk_word</i>	An integer containing a hook ID and a hook type: <table border="0" style="margin-left: 20px;"> <tr> <td><b>hk_id</b></td> <td>Before AIX 6.1 the hook identifier is a 12-bit value. On AIX 6.1 and above, the hook identifier is a 16-bit value. A 16-bit value of the form hhh0 is equivalent to a 12-bit value of the form hhh.</td> </tr> <tr> <td><b>hk_type</b></td> <td>A 4-bit hook type. The <b>trcgenk</b> service automatically records this information. This value is only valid before AIX 6.1.</td> </tr> </table>	<b>hk_id</b>	Before AIX 6.1 the hook identifier is a 12-bit value. On AIX 6.1 and above, the hook identifier is a 16-bit value. A 16-bit value of the form hhh0 is equivalent to a 12-bit value of the form hhh.	<b>hk_type</b>	A 4-bit hook type. The <b>trcgenk</b> service automatically records this information. This value is only valid before AIX 6.1.
<b>hk_id</b>	Before AIX 6.1 the hook identifier is a 12-bit value. On AIX 6.1 and above, the hook identifier is a 16-bit value. A 16-bit value of the form hhh0 is equivalent to a 12-bit value of the form hhh.				
<b>hk_type</b>	A 4-bit hook type. The <b>trcgenk</b> service automatically records this information. This value is only valid before AIX 6.1.				
<i>data_word</i>	Specifies a word of user-defined data.				
<i>len</i>	Specifies the length in bytes of the buffer specified by the <i>buf</i> parameter.				
<i>buf</i>	Points to a buffer of trace data. The maximum amount of trace data is 4096 bytes.				

## Description

The **trcgenk** kernel service records a trace event if a trace session is active for the specified trace channel. If a trace session is not active, the **trcgenk** kernel service simply returns. The **trcgenk** kernel service is located in pinned kernel memory.

The **trcgenk** kernel service is used to record a trace entry consisting of an *hk\_word* entry, a *data\_word* entry, a variable number of bytes of trace data, and, in AIX 5L™ Version 5.3 with the 5300-05 Technology Level and above, a time stamp.

## Execution Environment

The **trcgenk** kernel service can be called from either the process or interrupt environment.

## Return Values

The **trcgenk** kernel service has no return values.

### Related reference:

“trcgenkt Kernel Service”

### Related information:

trace subroutine

trcgen subroutine

RAS Kernel Services

## trcgenkt Kernel Service

### Purpose

Records a trace event, including a time stamp, for a generic trace channel.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/trchkid.h>
```

```
void trcgenkt (chan, hk_word, data_word, len, buf)
unsigned int chan, hk_word, data_word, len;
char * buf;
```

## Parameters

Item	Description
<i>chan</i>	Specifies the channel number for the trace session. This number is obtained from the <b>trcstart</b> subroutine.
<i>hk_word</i>	An integer containing a hook ID and a hook type:  <b>hk_id</b> Before AIX 6.1 the hook identifier is a 12-bit value. On AIX 6.1 and above, the hook identifier is a 16-bit value. A 16-bit value of the form hhh0 is equivalent to a 12-bit value of the form hhh.  <b>hk_type</b> A 4-bit hook type. The <b>trcgenkt</b> service automatically records this information. This value is only valid before AIX 6.1.
<i>data_word</i>	Specifies a word of user-defined data.
<i>len</i>	Specifies the length, in bytes, of the buffer identified by the <i>buf</i> parameter.
<i>buf</i>	Points to a buffer of trace data. The maximum amount of trace data is 4096 bytes.

## Description

The **trcgenkt** kernel service records a trace event if a trace session is active for the specified trace channel. If a trace session is not active, the **trcgenkt** service simply returns. The **trcgenkt** kernel service is located in pinned kernel memory.

The **trcgenkt** service records a trace entry consisting of an *hk\_word* entry, a *data\_word* entry, a variable number of bytes of trace data, and a time stamp.

## Execution Environment

The **trcgenkt** kernel service can be called from either the process or interrupt environment.

## Return Values

The **trcgenkt** service has no return values.

### Related reference:

“trcgenk Kernel Service” on page 493

### Related information:

trace command

trcgen subroutine

RAS Kernel Services

## **trcgenkt** Kernel Service for Data Link Control (DLC) Devices Purpose

Records a trace event, including a time stamp, for a DLC trace channel.

## Syntax

```
#include <sys/trchkid.h>
```

```
void trcgenkt (chan, hk_word, data_word, len, buf)  
unsigned int chan, hk_word, data_word, len;  
char * buf;
```

## Parameters

<b>Item</b>	<b>Description</b>
<i>chan</i>	Specifies the channel number for the trace session. This number is obtained from the <b>trcstart</b> subroutine.
<i>hk_word</i>	Contains the trace hook identifier defined in the <b>/usr/include/sys/trchkid.h</b> file. The types of link trace entries registered using the hook ID include: <ul style="list-style-type: none"> <li><b>HKWD_SYSX_DLC_START</b> Start link station completions</li> <li><b>HKWD_SYSX_DLC_TIMER</b> Time-out completions</li> <li><b>HKWD_SYSX_DLC_XMIT</b> Transmit completions</li> <li><b>HKWD_SYSX_DLC_RECV</b> Receive completions</li> <li><b>HKWD_SYSX_DLC_HALT</b> Halt link station completions</li> </ul>
<i>data_word</i>	Specifies trace data format field. This field varies depending on the hook ID. Each of these definitions are in the <b>/usr/include/sys/gdlextc.h</b> file: <ul style="list-style-type: none"> <li>• The first half-word always contains the data link protocol field including one of these definitions: <ul style="list-style-type: none"> <li><b>DLC_DL_SDLC</b> SDLC</li> <li><b>DLC_DL_HDLC</b> HDLC</li> <li><b>DLC_DL_BSC</b> BISYNC</li> <li><b>DLC_DL_ASC</b> ASYNC</li> <li><b>DLC_DL_PCNET</b> PC Network</li> <li><b>DLC_DL_ETHER</b> Standard Ethernet</li> <li><b>DLC_DL_802_3</b> IEEE 802.3</li> <li><b>DLC_DL_TOKEN</b> Token-Ring</li> </ul> </li> </ul>



**Item****Description**

- On start or halt link station completion, the second half-word contains the physical link protocol in use:

**DLC\_PL\_EIA232**

EIA-232D Telecommunications

**DLC\_PL\_EIA366**

EIA-366 Auto Dial

**DLC\_PL\_X21**

CCITT X.21 Data Network

**DLC\_PL\_PCNET**

PC Network Broadband

**DLC\_PL\_ETHER**

Standard Baseband Ethernet

**DLC\_PL\_SMART**

Smart Modem Auto Dial

**DLC\_PL\_802\_3**

IEEE 802.3 Baseband Ethernet

**DLC\_PL\_TBUS**

IEEE 802.4 Token Bus

**DLC\_PL\_TRING**

IEEE 802.5 Token-Ring

**DLC\_PL\_EIA422**

EIA-422 Telecommunications

**DLC\_PL\_V35**

CCITT V.35 Telecommunications

**DLC\_PL\_V25BIS**

CCITT V.25 bis Autodial for Telecommunications

- On timeout completion, the second half-word contains the type of timeout occurrence:

**DLC\_TO\_SLOW\_POLL**

Slow station poll

**DLC\_TO\_IDLE\_POLL**

Idle station poll

**DLC\_TO\_ABORT**

Link station aborted

**DLC\_TO\_INACT**

Link station receive inactivity

**DLC\_TO\_FAILSAFE**

Command failsafe

**DLC\_TO\_REPOLL\_T1**

Command repoll

**DLC\_TO\_ACK\_T2**

I-frame acknowledgment

- On transmit completion, the second half-word is set to the data link control bytes being sent. Some transmit packets only have a single control byte; in that case, the second control byte is not displayed.
  - On receive completion, the second half-word is set to the data link control bytes that were received. Some receive packets only have a single control byte; in that case, the second control byte is not displayed.
- Specifies the length in bytes of the entry specific data specified by the *buf* parameter.

*len*

Item	Description
<i>buf</i>	<p>Specifies the pointer to the entry specific data that consists of:</p> <p><b>Start Link Station Completions</b> Link station diagnostic tag and the remote station's name and address.</p> <p><b>Time-out Completions</b> No specific data is recorded.</p> <p><b>Transmit Completions</b> Either the first 80 bytes or all the transmitted data, depending on the short/long trace option.</p> <p><b>Receive Completions</b> Either the first 80 bytes or all the received data, depending on the short/long trace option.</p> <p><b>Halt Link Station Completions</b> Link station diagnostic tag, the remote station's name and address, and the result code.</p>

## Description

The **trcgenkt** kernel service records a trace event if a trace session is active for the specified trace channel. If a trace session is not active, the **trcgenkt** kernel service simply returns. The **trcgenkt** kernel service is located in pinned kernel memory.

The **trcgenkt** kernel service is used to record a trace entry consisting of an *hk\_word* entry, a *data\_word* entry, a variable number of bytes of trace data, and a time stamp.

## Execution Environment

The **trcgenkt** kernel service can be called from either the process or interrupt environment.

## Return Values

The **trcgenkt** kernel service has no return values.

### Related reference:

“trcgenk Kernel Service” on page 493

“trcgenkt Kernel Service” on page 494

### Related information:

trace subroutine

Generic Data Link Control (GDLC) Environment Overview

RAS Kernel Services

## tstart Kernel Service

### Purpose

Submits a timer request.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/timer.h>
```

```
void tstart ( t)
struct trb *t;
```

### Parameter

Item	Description
<i>t</i>	Points to a timer request structure.

## Description

The **tstart** kernel service submits a timer request with the timer request block specified by the *t* parameter as input. The caller of the **tstart** kernel service must first call the **talloc** kernel service to allocate the timer request structure. The caller must then initialize the structure's fields before calling the **tstart** kernel service.

Once the request has been submitted, the kernel calls the *t->func* timer function when the amount of time specified by the *t->timeout.it* value has elapsed. The *t->func* timer function is called on an interrupt level. Therefore, code for this routine must follow conventions for interrupt handlers.

The **tstart** kernel service examines the *t->flags* field to determine if the timer request being submitted represents an absolute request or an incremental one. An absolute request is a request for a time out at the time represented in the **it\_value** structure. An incremental request is a request for a time out at the time represented by now, plus the time in the **it\_value** structure.

The caller should place time information for both absolute and incremental timers in the **itimerstruc\_t.t.it** value substructure. The **T\_ABSOLUTE** absolute request flag is defined in the **/usr/include/sys/timer.h** file and should be ORed into the *t->flag* field if an absolute timer request is desired.

When the **T\_MOVE\_OK** flag is set, the associated timer is moved to another processor when the owning processor is folded.

When **T\_LATE\_OK** flag is set, the associated timer is put to sleep when the owning processor is put to sleep (folded) mode. The timer expiration handler is called when the owning processor is awakened (unfolded) if the scheduled expiration time has past. The time spent sleeping is therefore counted with respect to the expiration time. When this flag is set, there is no guarantee as to when the timer might expire.

**Note:** The **T\_MOVE\_OK** and **T\_LATE\_OK** flags are not required. They are intended to improve the effectiveness of processor folding by reducing the load on folded processors.

Modifications to the system time are added to incremental timer requests, but not to absolute ones. Consider the user who has submitted an absolute timer request for noon on 12/25/88. If a privileged user then modifies the system time by adding four hours to it, then the timer request submitted by the user still occurs at noon on 12/25/88.

By contrast, suppose it is presently 12 noon and a user submits an incremental timer request for 6 hours from now (to occur at 6 p.m.). If, before the timer expires, the privileged user modifies the system time by adding four hours to it, the user's timer request will then expire at 2200 (10 p.m.).

## Execution Environment

The **tstart** kernel service can be called from either the process or interrupt environment.

## Return Values

The **tstart** service has no return values.

### Related reference:

“tstop Kernel Service” on page 500

### Related information:

## **tstop Kernel Service**

### **Purpose**

Cancels a pending timer request.

### **Syntax**

```
#include <sys/types.h>  
#include <sys/errno.h>  
#include <sys/timer.h>
```

```
int tstop ( t )  
struct trb *t;
```

### **Parameter**

Item	Description
<i>t</i>	Specifies the pending timer request to cancel.

### **Description**

The **tstop** kernel service cancels a pending timer request. The **tstop** kernel service must be called before a timer request block can be freed with the **tfree** kernel service.

In a multiprocessor environment, the timer function associated with a timer request block may be active on another processor when the **tstop** kernel service is called. In this case, the timer request cannot be canceled. A multiprocessor-safe driver must therefore check the return code and take appropriate action if the cancel request failed.

In a uniprocessor environment, the call always succeeds. This is untrue in a multiprocessor environment, where the call will fail if the timer is being handled by another processor. Therefore, the function now has a return value, which is set to 0 if successful, or -1 otherwise. Funnelled device drivers do not need to check the return value since they run in a logical uniprocessor environment. Multiprocessor-safe and multiprocessor-efficient device drivers need to check the return value in a loop. In addition, if a driver uses locking, it must release and reacquire its lock within this loop. A delay should be used between the release and reacquiring the lock as shown below:

```
while (tstop(&trp)) {  
    release_any_lock;  
    delay_some_time;  
    reacquire_the_lock;  
} /* null while loop if locks not used */
```

### **Execution Environment**

The **tstop** kernel service can be called from either the process or interrupt environment.

### **Return Values**

Item	Description
0	Indicates that the request was successfully canceled.
-1	Indicates that the request could not be canceled.

#### Related reference:

“tstart Kernel Service” on page 498

#### Related information:

Timer and Time-of-Day Kernel Services

Using Fine Granularity Timer Services and Structures

Using Multiprocessor-Safe Timer Services

## tuning Kernel Service

### Purpose

Provides access to the kernel tunable variables through an easily accessible interface.

### Syntax

```
typedef enum {
    TH_MORE,
    TH_EOF
} tmode_t;

#define TH_ABORT TH_EOF

typedef int (*tuning_read_t)(tmode_t mode, long *size, char **buf, void *context);
typedef int (*tuning_write_t)(tmode_t mode, long *size, char *buf, void *context);

tinode_t *tuning_register_handler (path, mode, readfunc, writefunc, context)
const char *path;
mode_t mode;
tuning_read_t readfunc;
tuning_write_t writefunc;
void * context;

tinode *tuning_register_bint32 (path, mode, variable, low, high)
const char *path;
mode_t mode;
int32 *variable;
int32 low;
int32 high;

tinode *tuning_register_bint32x (path, rfunc, wfunc, mode, low, high)
const char *path;
mode_t mode;
int32 (*rfunc)(void *);
int (*wfunc)(int32, void *);
void *context;
int32 low;
int32 high;

tinode *tuning_register_buint32 (path, mode, variable, low, high)
const char *path;
mode_t mode;
uint32 *variable;
uint32 low;
uint32 high;

tinode *tuning_register_buint32x (path, rfunc, wfunc, mode, low, high)
const char *path;
mode_t mode;
uint32 (*rfunc)(void *);
int (*wfunc)(uint32, void *);
void *context;
uint32 low;
uint32 high;
```

```

tinode *tuning_register_bint64 (path, mode, variable, low, high)
const char *path;
mode_t mode;
int64 *variable;
int64 low;
int64 high;

tinode *tuning_register_bint64x (path, rfunc, wfunc, mode, low, high)
const char *path;
mode_t mode;
int64 (*rfunc)(void *);
int (*wfunc)(int64, void *);
void *context;
int64 low;
int64 high;

tinode *tuning_register_buint64 (path, mode, variable, low, high)
const char *path;
mode_t mode;
uint64 *variable;
uint64 low;
uint64 high;

tinode *tuning_register_buint64x (path, rfunc, wfunc, mode, low, high)
const char *path;
mode_t mode;
uint64 (*rfunc)(void *);
int (*wfunc)(uint64, void *);
void *context;
uint64 low;
uint64 high;

void tuning_deregister (t)
tinode_t * t;

```

## Description

The **tuning\_register\_handler** kernel service is used to add a file at the location specified by the *path* parameter. When this file is read from or written to, one of the two callbacks passed as parameters to the function is invoked.

Accesses to the file are viewed in terms of streams. A single stream is created by a sequence of one open, one or more reads, and one close on the file. While the file is open by one process, attempts to open the same file by other processes will be blocked unless **O\_NONBLOCK** is passed in the flags to the **open** subroutine.

The *readfunc* callback behaves like a producer function. The function is called when the user attempts to read from the file. The *mode* parameter is equal to **TH\_MORE** unless the user closes the file prematurely. On entry, the *size* parameter is an integer containing the size of the buffer. The *context* parameter is the context pointer passed to the registration function. Upon return, *size* should contain either the actual amount of data returned, or a zero if an end-of-file condition should be returned to the user. The return value of the function can also be used to signal end-of-file, as described below.

**Note:** It is expected that the *readfunc* callback has already done any necessary end-of-file cleanup when it returns the end-of-file signal.

If the amount of data returned is nonzero, the *buf* parameter may be modified to point to a new buffer. If this is done, the callback is responsible for freeing the new buffer.

If the buffer provided by the caller is too small, the caller may instead set *buf* to NULL. In this case, the *size* parameter should be modified to indicate the size of the buffer needed. The caller will then re-invoke the callback with a buffer of at least the requested size.

If the user closes the file before the callback indicates end-of-file, the callback will be invoked one last time with *mode* equal to **TH\_ABORT**. In this case, the *size* parameter is equal to 0 on entry, and any data returned is discarded. The callback must reset its state because no further callbacks will be made for this stream.

The *writefunc* callback behaves as a consumer function and is used when the user attempts to write to the file. The *mode* parameter is set to **TH\_EOF** if no further data can be expected on this stream (for example, the user called the **close** subroutine on the file). Otherwise, *mode* is set to **TH\_MORE**. The *size* parameter contains the size of the data passed in the buffer. The *buf* parameter is the pointer to the buffer.

**Note:** There will be zero or more calls with the *mode* parameter set to **TH\_MORE** and one call with the *mode* parameter set to **TH\_EOF** for every stream.

The *buf* parameter may change between invocations. Upon return from the callback, the *size* parameter must be modified to reflect the amount of data consumed from the buffer, and the buffer must not be freed even if all data is consumed. The function is expected to consume data in a linear (first in, first out) fashion. Unconsumed data is present at the beginning of the buffer at the next invocation of the callback. The *size* parameter will include the size of the unconsumed data.

Both callbacks' return values are expected to be zero. If unsuccessful, a positive value will be placed into the **errno** global variable (with the accompanying indication of an error return from the kernel service). If the return value of a callback is less than 0, end-of-file will be signaled to the user, and the return value will be treated as its unary negation (For example, -1 will be treated like 0). In this case, no further callbacks will be made for this stream.

The **tuning\_register\_bint32**, **tuning\_register\_buint32**, **tuning\_register\_bint64**, and **tuning\_register\_buint64** kernel services are used to add a file at the location specified by the *path* parameter that, when read from, will return the ASCII value of the integer variable pointed to by the *variable* parameter. When written to, this file will set the integer variable to the value whose ASCII value was written, unless that value does not satisfy the relation  $low \leq value < high$ . In this case, the integer variable is not modified, and an error is returned to the user through an error return of the kernel service during which the invalid attempt is detected (probably either **write** or **close**).

The **tuning\_register\_b\*x** functions operate similarly to their non-*x* variants, but they use a pair of callbacks to retrieve (*rfunc*) and set (*wfunc*) the variable. The callback is passed the value (if setting) and the context parameter. This permits more complex operations on read/write, such as serialization and memory allocation and deallocation.

The **tuning\_get\_context** kernel service returns the *context* of the registration function used to create the **tinode\_t** structure referred to by the *argument* parameter.

The **tuning\_register** kernel service is the basic interface by which a file can be added to the **/proc/sys** directory hierarchy. This function is not exported to kernel extensions, and its direct use in the kernel is strongly discouraged. The *path* parameter contains the path relative to the **/proc/sys** root at which the file should appear. Intermediate path components are automatically created. The *mode* parameter contains the UNIX permissions and the type of the file to be created (as per the **st\_mode** field of the **stat** struct). If the file type is not specified, it is assumed to be **S\_IFREG**. In most cases this parameter will be 0644 or 0600. The *vnops* parameter is used to dispatch all operations on the file.

The **tuning\_deregister** kernel service is used to remove a file from the **/proc/sys** directory hierarchy. It is exported to kernel extensions. It should only be used when a specific file's implementation is no longer available. The *t* parameter is a **tinode\_t** structure as returned by **tuning\_register**. If the file is currently open, any further access to it after this call returns **ESTALE**.

## Parameters

Item	Description
<i>mode</i>	Is set to either <b>TH_EOF</b> if no further data is expected from the user for this change, or <b>TH_MORE</b> if further data is expected.
<i>size</i>	Contains the size of the data passed in the buffer.
<i>buf</i>	Points to the buffer.
<i>context</i>	Points to the context passed to the registration function.
<i>path</i>	Specifies the location of the file to be added.
<i>readfunc</i>	Behaves as a producer function.
<i>rfunc</i>	Retrieves the variable.
<i>wfunc</i>	Sets the variable.
<i>writefunc</i>	Behaves as a consumer function.
<i>variable</i>	Specifies the variable.
<i>high</i>	Specifies the maximum value that the <i>variable</i> parameter can contain.
<i>low</i>	Specifies the minimum value that the <i>variable</i> parameter can contain.
<i>t</i>	A <b>tinode_t</b> structure as returned by <b>tuning_register</b> .

## Return Values

Upon successful completion, the **tuning\_register** kernel service returns the newly created **tinode\_t** structure. If unsuccessful, a NULL value is returned.

## Examples

A user of this interface might include the following line in their initialization routine:

```
tuning_var = tuning_register_buint64
("fs/jfs2/max_readahead", 0644 &j2_max_read_ahead, 0, 1024);
```

In this example *tuning\_var* is a global variable of type **tinode\_t** \*. This causes the **fs** and **fs/jfs2** directories to be created, and a file (pipe) to be created as **fs/jfs2/max\_readahead**. The file returns the value of **j2\_max\_readahead** in ASCII when read. The variable is read at the time of the first read. A write would set the value of the variable, but only at the time of either the first newline being written or a **close** function being performed. In order to write the variable after reading it, one must close the file and reopen it for write. This file is not seekable.

## U

The following kernel services begin with the with the letter u.

### ue\_proc\_check Kernel Service

#### Purpose

Determines if a process is critical to the system.

#### Syntax

```
int ue_proc_check (pid)
pid_t pid;
```

#### Description

The **ue\_proc\_check** kernel service determines if a particular process is critical to the system. A critical process is either a kernel process or a process registered as critical by the **ue\_proc\_register** system call. A process that is critical will cause the system to terminate if that process has an unrecoverable hardware error associated with the process. Unrecoverable hardware errors associated with a process are determined by the kernel machine check handler on systems that support UE-Gard error processing.

The **ue\_proc\_check** kernel service should be called only while executing in kernel mode in the user process.



## Parameters

Item	Description
<i>pid</i>	Specifies the process' ID to be checked as critical.

## Execution Environment

The `ue_proc_check` kernel service can be called from the interrupt environment only.

## Return Values

Item	Description
0	Indicates that the <i>pid</i> is not critical.
EINVAL	Indicates that the <i>pid</i> is critical.
-1	Indicates that the <i>pid</i> parameter is not valid or the process no longer exists.

### Related reference:

“ue\_proc\_register Subroutine”

## ue\_proc\_register Subroutine

### Purpose

Registers a process as critical to the system.

### Syntax

```
int ue_proc_register (pid, argument)
pid_t pid;
int argument;
```

### Description

The `ue_proc_register` system call registers a particular process as critical to the system. A process that is critical will cause the system to terminate if that process has an unrecoverable hardware error associated with the process. Unrecoverable hardware errors associated with a process are determined by the kernel machine check handler on systems that support UE-Gard error processing.

An execed process from a critical process must register itself to be critical. A fork from a process inherits the critical registration unless the argument is set to **NONCRITFORK**.

If the value of the *pid* parameter is equal to (`pid_t`) 0, the subroutine is registering the calling process.

The `ue_proc_register` system call should be called only while executing with root authority in the user process.

## Parameters

Item	Description
<i>pid</i>	Specifies the process' ID to be registered critical.
<i>argument</i>	Defined in the <code>sys/proc.h</code> header file. Can be the following value: <b>NONCRITFORK</b> The <i>pid</i> forks are not critical.

## Execution Environment

The `ue_proc_register` system call can be called from the process environment only.

## Return Values

Item	Description
0	Indicates successful completion.
EINVAL	Indicates that the <i>pid</i> parameter is not valid or the process no longer exists.
EACCES	Indicates that the caller does not have sufficient authority to alter the <i>pid</i> registration.

### Related reference:

“ue\_proc\_unregister Subroutine”

## ue\_proc\_unregister Subroutine

### Purpose

Unregisters a process from being critical to the system.

### Syntax

```
int ue_proc_unregister (pid)
pid_t pid;
```

### Description

The **ue\_proc\_unregister** system call unregisters a particular process as being no longer critical to the system. A process that has been previously registered critical will cause the system to terminate if that process has an unrecoverable hardware error associated with the process. Unrecoverable hardware errors associated with a process are determined by the kernel machine check handler on systems that support UE-Gard error processing.

If the value of the *pid* parameter is equal to (**pid\_t**) 0, the subroutine is unregistering the calling process.

The **ue\_proc\_unregister** service should be called only while executing with root authority in the user process.

### Parameters

Item	Description
<i>pid</i>	Specifies the process' ID to be unregistered.

## Execution Environment

The **ue\_proc\_unregister** system call can be called from the process environment only.

## Return Values

Item	Description
0	Indicates successful completion.
EINVAL	Indicates that the <i>pid</i> parameter is not valid or the process no longer exists.
EACCES	Indicates that the caller does not have sufficient authority to alter the <i>pid</i> registration.

### Related reference:

“ue\_proc\_register Subroutine” on page 505

## uexadd Kernel Service

### Purpose

Adds a systemwide exception handler for catching user-mode process exceptions.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/except.h>
```

```
void uexadd ( exp)
struct uexcept *exp;
```

## Parameter

Item	Description
<i>exp</i>	Points to an exception handler structure. This structure must be pinned and is used for registering user-mode process exception handlers. The <b>uexcept</b> structure is defined in the <code>/usr/include/sys/except.h</code> file.

## Description

The **uexadd** kernel service is typically used to install a systemwide exception handler to catch exceptions occurring during execution of a process in user mode. The **uexadd** kernel service adds the exception handler structure specified by the *exp* parameter, to the chain of exception handlers to be called if an exception occurs while a process is executing in user mode. The last exception handler registered is the first exception handler called for a user-mode exception.

The **uexcept** structure has:

- A chain element used by the kernel to chain the registered user exception handlers.
- A function pointer defining the entry point of the exception handler being added.

Additional exception handler-dependent information can be added to the end of the structure, but must be pinned.

**Attention:** The **uexcept** structure must be pinned when the **uexadd** kernel service is called. It must remain pinned and unmodified until after the call to the **uexdel** kernel service to delete the specified exception handler. Otherwise, the system may crash.

## Execution Environment

The **uexadd** kernel service can be called from the process environment only.

## Return Values

The **uexadd** kernel service has no return values.

### Related reference:

“uexdel Kernel Service” on page 510

“User-Mode Exception Handler for the uexadd Kernel Service”

### Related information:

User-Mode Exception Handling

Kernel Extension and Device Driver Management Services

## User-Mode Exception Handler for the uexadd Kernel Service Purpose

Handles exceptions that occur while a kernel thread is executing in user mode.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/except.h>
```

```
int func (exp, type, tid, mst)
struct except * exp;
int type;
tid_t tid;
struct kmstsave * mst;
```

## Parameters

Item	Description
------	-------------

<i>exp</i>	Points to the <b>except</b> structure used to register this exception handler.
<i>mst</i>	Points to the current <b>kmstsave</b> area for the process. This pointer can be used to access the <b>kmstsave</b> area to obtain additional information about the exception.

Item	Description
------	-------------

<i>tid</i>	Specifies the thread ID of the kernel thread that was executing at the time of the exception.
<i>type</i>	Denotes the type of exception that has occurred. This type value is platform specific. Specific values are defined in the <code>/usr/include/sys/except.h</code> file.

## Description

The user-mode exception handler (*exp*->**func**) is called for synchronous exceptions that are detected while a kernel thread is executing in user mode. The kernel exception handler saves exception information in the **kmstsave** area of the structure. For user-mode exceptions, it calls the first exception handler found on the user exception handler list. The exception handler executes in an interrupt environment at the priority level of either **INTPAGER** or **INTIODONE**.

If the registered exception handler returns a return code indicating that the exception was handled, the kernel exits from the exception handler without calling additional exception handlers from the list. If the exception handler returns a return code indicating that the exception was not handled, the kernel invokes the next exception handler on the list. The last exception handler in the list is the default handler. This is typically signalling the thread.

The kernel exception handler must not page fault. It should also register an exception handler using the **setjmpx** kernel service if any exception-handling activity can result in an exception. This is important particularly if the exception handler is handling the I/O. If the exception handler did not handle the exception, the return code should be set to the **EXCEPT\_NOT\_HANDLED** value for user-mode exception handling.

## Execution Environment

The user-mode exception handler for the **uexadd** kernel service is called in the interrupt environment at the **INTPAGER** or **INTIODONE** priority level.

## Return Values

Item	Description
EXCEPT_HANDLED	Indicates that the exception was successfully handled.
EXCEPT_NOT_HANDLED	Indicates that the exception was not handled.

**Related reference:**

“uexadd Kernel Service” on page 506

**Related information:**

User-Mode Exception Handling

Kernel Extension and Device Driver Management Kernel Services

## uexblock Kernel Service

### Purpose

Makes the currently active kernel thread nonrunnable when called from a user-mode exception handler.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/except.h>
```

```
void uexblock ( tid)
tid_t *tid;
```

### Parameter

Item	Description
<i>tid</i>	Specifies the thread ID of the currently active kernel thread to be put into a wait state.

### Description

The **uexblock** kernel service puts the currently active kernel thread specified by the *tid* parameter into a wait state until the **uexclear** kernel service is used to make the thread runnable again. If the **uexblock** kernel service is called from the process environment, the *tid* parameter must specify the current active thread; otherwise the system will crash with a kernel panic.

The **uexblock** kernel service can be used to lazily control user-mode threads access to a shared serially usable resource. Multiple threads can use a serially used resource, but only one process at a time. When a thread attempts to but cannot access the resource, a user-mode exception can be set up to occur. This gives control to an exception handler registered by the **uexadd** kernel service. This exception handler can then block the thread using the **uexblock** kernel service until the resource is made available. At this time, the **uexclear** kernel service can be used to make the blocked thread runnable.

### Execution Environment

The **uexblock** kernel service can be called from either the process or interrupt environment.

### Return Values

The **uexblock** service has no return values.

**Related reference:**

“uexclear Kernel Service” on page 510

**Related information:**

User-Mode Exception Handling

Kernel Extension and Device Driver Management Services

## uexclear Kernel Service

### Purpose

Makes a kernel thread blocked by the **uexblock** service runnable again.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/except.h>
```

```
void uexclear ( tid)
tid_t *tid;
```

### Parameter

Item	Description
<i>tid</i>	Specifies the thread ID of the previously blocked kernel thread to be put into a run state.

### Description

The **uexclear** kernel service puts a kernel thread specified by the *tid* parameter back into a runnable state after it was made nonrunnable by the **uexblock** kernel service. A thread that has been sent a **SIGSTOP** stop signal is made runnable again when it receives the **SIGCONT** continuation signal.

The **uexclear** kernel service can be used to lazily control user-mode thread access to a shared serially usable resource. A serially used resource is usable by more than one thread, but only by one at a time. When a thread attempts to access the resource but does not have access, a user-mode exception can be setup to occur.

This setup gives control to an exception handler registered by the **uexadd** kernel service. Using the **uexblock** kernel service, this exception handler can then block the thread until the resource is later made available. At that time, the **uexclear** service can be used to make the blocked thread runnable.

### Execution Environment

The **uexclear** kernel service can be called from either the process or interrupt environment.

### Return Values

The **uexclear** service has no return values.

#### Related reference:

“uexblock Kernel Service” on page 509

#### Related information:

User-Mode Exception Handling

Kernel Extension and Device Driver Management Services

## uexdel Kernel Service

### Purpose

Deletes a previously added systemwide user-mode exception handler.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/except.h>
```

```
void uexdel ( exp)
struct uexcepth *exp;
```

## Parameter

Item	Description
------	-------------

<i>exp</i>	Points to the exception handler structure used to add the exception handler with the <b>uexadd</b> kernel service.
------------	--

## Description

The **uexdel** kernel service removes a user-mode exception handler from the systemwide list of exception handlers maintained by the kernel's exception handler.

The **uexdel** kernel service removes the exception handler structure specified by the *exp* parameter from the chain of exception handlers to be called if an exception occurs while a process is executing in user mode. Once the **uexdel** kernel service has completed, the specified exception handler is no longer called. In addition, the **uexcepth** structure can be modified, freed, or unpinned.

## Execution Environment

The **uexdel** kernel service can be called from the process environment only.

## Return Values

The **uexdel** kernel service has no return values.

### Related reference:

“uexadd Kernel Service” on page 506

### Related information:

User-Mode Exception Handling

Kernel Extension and Device Driver Management Services

## ufdcreate Kernel Service

### Purpose

Allocates and initializes a file descriptor.

### Syntax

```
#include <fcntl.h>
#include <sys/types.h>
#include <sys/file.h>
int ufdcreate (flags, ops, datap, type, fdp, cnp)
```

```
int flags;
struct fileops * ops;
void * datap;
short type;
int * fdp;
struct ucred *crp;
```

### Parameters

Item	Description
<i>flags</i>	Specifies the flags to save in a <b>file</b> structure. The <b>file</b> structure is defined in the <b>sys/file.h</b> file. If a <b>read</b> or <b>write</b> subroutine is called with the file descriptor returned by this routine, the <b>FREAD</b> and <b>FWRITE</b> flags must be set appropriately. Valid flags are defined in the <b>fcntl.h</b> file.
<i>ops</i>	Points to the list of subsystem-supplied routines to call for the file system operations: read/write, ioctl, select, fstat, and close. The <b>fileops</b> structure is defined in the <b>sys/file.h</b> file. See "File Operations" for more information.
<i>datap</i>	Points to type-dependent structures. The system saves this pointer in the <b>file</b> structure. As a result, the pointer is available to the file operations when they are called.
<i>type</i>	Specifies the unique type value for the <b>file</b> structure. Valid types are listed in the <b>sys/file.h</b> file.
<i>fdp</i>	Points to an integer field where the file descriptor is stored on successful return.
<i>crp</i>	Points to a credentials structure. This pointer is saved in the file struct for use in subsequent operations. It must be a valid <b>ucred</b> struct. The <b>crref()</b> kernel service can be used to obtain a <b>ucred</b> struct.

## Description

The **ufdcreate** kernel service provides a file interface to kernel extensions. Kernel extensions use this service to create a file descriptor and file structure pair. Also, this service allows kernel extensions to provide their own file descriptor-based system calls, enabling read/write, ioctl, select, fstat, and close operations on objects outside the file system. The **ufdcreate** kernel services does not require the extension to understand or conform to the synchronization requirements of the logical file system (LFS).

The **ufdcreate** kernel service provides a file descriptor to the caller and creates the underlying file structure. The caller must include pointers to subsystem-supplied routines for the read/write, ioctl, select, fstat, and close operations. If any of the operations are not needed by the calling subsystem, then the caller must provide a pointer to an appropriate **errno** value. Typically, the **EOPNOTSUPP** value is used for this purpose. See "File Operations" for information about the requirements for the subsystem-supplied routines.

## Removing a File Descriptor

There is no corresponding operation to remove a file descriptor (and the attendant structures) created by the **ufdcreate** kernel service. To remove a file descriptor, use a call to the **close** subroutine. The **close** subroutine can be called from a routine or from within the kernel or kernel extension. If the close is not called, the file is closed when the process exits.

Once a call is made to the **ufdcreate** kernel service, the file descriptor is considered open before the call to the service returns. When a **close** or **exit** subroutine is called, the close file operation specified on the call to the **ufdcreate** interface is called.

## File Operations

The **ufdcreate** kernel service allows kernel extensions to provide their own file descriptor-based system calls, enabling read/write, ioctl, select, fstat, and close operations on objects outside the file system. The **fileops** structure defined in the **sys/file.h** file provides interfaces for these routines.

### read/write Requirements

The read/write operation manages input and output to the object specified by the *fp* parameter. The actions taken by this operation are dependent on the object type. The syntax for the operation is as follows:

```
#include <sys/types.h>
#include <sys/uio.h>
int (*fo_rw) (fp, rw, uiop, ext)
```



```

struct file *fp;
enum uio_rw rw;
struct uio *uiop;
int ext;

```

The parameters have the following values:

Value	Description
<i>fp</i>	Points to the <b>file</b> structure. This structure corresponds to the file descriptor used on the <b>read</b> or <b>write</b> subroutine.
<i>rw</i>	Contains a <b>UIO_READ</b> value for a read operation or <b>UIO_WRITE</b> value for a write operation.
<i>uiop</i>	Points to a <b>uio</b> structure. This structure describes the location and size information for the input and output requested. The <b>uio</b> structure is defined in the <b>uio.h</b> file.
<i>ext</i>	Specifies subsystem-dependent information. If the <b>readx</b> or <b>writex</b> subroutine is used, the value passed by the operation is passed through to this subroutine. Otherwise, the value is 0.

If successful, the **fo\_rw** operation returns a value of 0. A nonzero return value should be programmed to indicate an error. See the **sys/errno.h** file for a list of possible values.

**Note:** On successful return, the `uiop->uio_resid` field must be updated to include the number of bytes of data actually transferred.

## ioctl Requirements

The **ioctl** operation provides object-dependent special command processing. The **ioctl** subroutine performs a variety of control operations on the object associated with the specified open **file** structure. This subroutine is typically used with character or block special files and returns an error for ordinary files.

The control operation provided by the **ioctl** operation is specific to the object being addressed, as are the data type and contents of the *arg* parameter.

The syntax for the **ioctl** operation is as follows:

```

#include <sys/types.h>
#include <sys/ioctl.h>
int (*fo_ioctl) (fp, cmd, arg, ext, kflag)
struct file *fp;
int cmd, ext, kflag;
caddr_t arg;

```

The parameters have the following values:

Value	Description
<i>fp</i>	Points to the <b>file</b> structure. This structure corresponds to the file descriptor used by the <b>ioctl</b> subroutine.
<i>cmd</i>	Defines the specific request to be acted upon by this routine.
<i>arg</i>	Contains data that is dependent on the <i>cmd</i> parameter.
<i>ext</i>	Specifies subsystem-specific information. If the <b>ioctlx</b> subroutine is used, the value passed by the application is passed through to this subroutine. Otherwise, the value is 0.
<i>kflag</i>	Determines where the call is made from. The <i>kflag</i> parameter has the value <b>FKERNEL</b> (from the <b>fcntl.h</b> file) if this routine is called through the <b>fp_ioctl</b> interface. Otherwise, its value is 0.

If successful, the **fo\_ioctl** operation returns a value of 0. For errors, the **fo\_ioctl** operation should return a nonzero return value to indicate an error. Refer to the **sys/errno.h** file for the list of possible values.

## select Requirements

The **select** operation performs a **select** operation on the object specified by the *fp* parameter. The syntax for this operation is as follows:

```

#include <sys/types.h>
int (*fo_select) (fp, corl, reqevents, rtneventsp, notify)
struct file *fp;
int corl;
ushort reqevents, *rtneventsp;
void (notify) ();

```

The parameters have the following values:

Value	Description
<i>fp</i>	Points to the <b>file</b> structure. This structure corresponds to the file descriptor used by the <b>select</b> subroutine.
<i>corl</i>	Specifies the ID used for correlation in the <b>selnotify</b> kernel service.
<i>reqevents</i>	Identifies the events to check. The poll and select functions define three standard event flags and one informational flag. The <b>sys/poll.h</b> file details the event bit definition. See the <b>fp_select</b> kernel service for information about the possible flags.
<i>rtneventsp</i>	Indicates the returned events pointer. This parameter, passed by reference, indicates the events that are true at the current time. The returned event bits include the request events and an error event indicator.
<i>notify</i>	Points to a routine to call when the specified object invokes the <b>selnotify</b> kernel service for an outstanding asynchronous select or poll event request. If no routine is to be called, this parameter must be null.

If successful, the **fo\_select** operation returns a value of 0. This operation should return a nonzero return value to indicate an error. Refer to the **sys/errno.h** file for the list of possible values.

### fstat Requirements

The **fstat** operation fills in an **attribute** structure. Depending on the object type specified by the *fp* parameter, many fields in the structure may not be applicable. The value passed back from this operation is dependent upon both the object type and what any routine that understands the type is expecting. The syntax for this operation is as follows:

```

#include <sys/types.h>
int (*fo_fstat) (fp, sbp)
struct file *fp;
struct stat *sbp;

```

The parameters have the following values:

Value	Description
<i>fp</i>	Points to the <b>file</b> structure. This structure corresponds to the file descriptor used by the <b>stat</b> subroutine.
<i>sbp</i>	Points to the <b>stat</b> structure to be filled in by this operation. The address supplied is in kernel space.

If successful, the **fo\_fstat** operation returns a value of 0. A nonzero return value should be programmed to indicate an error. Refer to the **sys/errno.h** file for the list of possible values.

### close Requirements

The **close** operation invalidates routine access to objects specified by the *fp* parameter and releases any data associated with that access. This operation is called from the **close** subroutine code when the **file** structure use count is decremented to 0. For example, if there are multiple accesses to an object (created by the **dup**, **fork**, or other subsystem-specific operation), the **close** subroutine calls the **close** operation when it determines that there is no remaining access through the **file** structure being closed.

A file descriptor is considered open once a file descriptor and **file** structure have been set up by the LFS. The **close** file operation is called whenever a **close** or **exit** is specified. As a result, the **close** operation must be able to close an object that is not fully open, depending on what the caller did before the **file** structure was initialized.

The syntax for the close operation is as follows:

```
#include <sys/file.h>
int (*fo_close) (fp)
struct file *fp;
```

The parameter is:

Item	Description
<i>fp</i>	Points to the <b>file</b> structure. This structure corresponds to the file descriptor used by the <b>close</b> subroutine.

If successful, the **fo\_close** operation returns a value of 0. This operation should return a nonzero return value to indicate an error. Refer to the **sys/errno.h** file for the list of possible values.

## Execution Environment

The **ufdcreate** kernel service can be called from the process environment only.

## Return Values

If the **ufdcreate** kernel service succeeds, it returns a value of 0. If the kernel service fails, it returns a nonzero value and sets the **errno** global variable.

## Error Codes

The **ufdcreate** kernel service fails if one or more of the following errors occur:

Error	Description
EINVAL	The <i>ops</i> parameter is null, or the <b>fileops</b> structure does not have entries for every operation.
EMFILE	All file descriptors for the process have already been allocated.
ENFILE	The system file table is full.

## Related reference:

“selnotify Kernel Service” on page 463

## Related information:

close subroutine

exit, atexit, or \_exit

Logical File System Kernel Services

## ufdgetf Kernel Service

### Purpose

Returns a pointer to a file structure associated with a file descriptor.

### Syntax

```
#include <sys/file.h>

int ufdgetf( fd, fpp)
int fd;
struct file **fpp;
```

### Parameters

Item	Description
<i>fd</i>	Identifies the file descriptor. The descriptor must be for an open file.
<i>fpp</i>	Points to a location to store the file pointer.

## Description

The **ufdgetf** kernel service returns a pointer to a file structure associated with a file descriptor. The calling routine must have a use count on the file descriptor. To obtain a use count on the file descriptor, the caller must first call the **ufdhold** kernel service.

## Execution Environment

The **ufdget** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates successful completion.
EBADF	Indicates that the <i>fd</i> parameter is not a file descriptor for an open file.

### Related reference:

“ufdhold and ufdrele Kernel Service”

## ufdhold and ufdrele Kernel Service

### Purpose

Increment or decrement a file descriptor reference count.

### Syntax

```
int ufdhold( fd)
int fd;
int ufdrele(fd)
int fd;
```

### Parameter

Item	Description
<i>fd</i>	Identifies the file descriptor.

## Description

**Attention:** It is extremely important that the calls to **ufdhold** and **ufdrele** kernel service are balanced. If a file descriptor is held more times than it is released, the **close** subroutine on the descriptor never completes. The process hangs and cannot be killed. If the descriptor is released more times than it is held, the system panics.

The **ufdhold** and **ufdrele** kernel services increment and decrement a file-descriptor reference count. Together, these kernel services maintain the file descriptor reference count. The **ufdhold** kernel service increments the count. The **ufdrele** kernel service decrements the count.

These subroutines are supported for kernel extensions that provide their own file-descriptor-based system calls. This support is required for synchronization with the **close** subroutine.

When a thread is executing a file-descriptor-based system call, it is necessary that the logical file system (LFS) be aware of it. The LFS uses the count in the file descriptor to monitor the number of system calls currently using any particular file descriptor. To keep the count accurately, any thread using the file descriptor must increment the count before performing any operation and decrement the count when all activity using the file descriptor is completed for that system call.

## Execution Environment

These kernel services can be called from the process environment only.

## Return Values

Item	Description
0	Indicates successful completion.
EBADF	Indicates that the <i>fd</i> parameter is not a file descriptor for an open file.

### Related reference:

“*ufdgetf* Kernel Service” on page 515

### Related information:

*close* subroutine

## uiomove Kernel Service

### Purpose

Moves a block of data between kernel space and a space defined by a **uio** structure.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/uio.h>
```

```
int uiomove ( cp, n, rw, uiop)
caddr_t cp;
int n;
uio_rw rw;
struct uio *uiop;
```

### Parameters

Item	Description
<i>cp</i>	Specifies the address in kernel memory to or from which data is moved.
<i>n</i>	Specifies the number of bytes to move.
<i>rw</i>	Indicates the direction of the move: <b>UIO_READ</b> Copies data from kernel space to space described by the <b>uio</b> structure. <b>UIO_WRITE</b> Copies data from space described by the <b>uio</b> structure to kernel space.
<i>uiop</i>	Points to a <b>uio</b> structure describing the buffer used in the data transfer.

### Description

The **uiomove** kernel service moves the specified number of bytes of data between kernel space and a space described by a **uio** structure. Device driver top halves, especially character device drivers,

frequently use the **uiomove** service to transfer data into or out of a user area. The `uio_resid` and `uio_iovcnt` fields in the **uio** structure describing the data area must be greater than 0 or an error is returned.

The **uiomove** service moves the number of bytes of data specified by either the `n` or `uio_resid` parameter, whichever is less. If either the `n` or `uio_resid` parameter is 0, no data is moved. The `uio_segflg` field in the **uio** structure is used to indicate if the move is accessing a user- or kernel-data area, or if the caller requires cross-memory operations and has provided the required cross-memory descriptors. If a cross-memory operation is indicated, there must be a cross-memory descriptor in the `uio_xmem` array for each `iovec` element.

If the move is successful, the following fields in the **uio** structure are updated:

Field	Description
<code>uio_iov</code>	Specifies the address of current <code>iovec</code> element to use.
<code>uio_xmem</code>	Specifies the address of the current <code>xmem</code> element to use.
<code>uio_iovcnt</code>	Specifies the number of remaining <code>iovec</code> elements.
<code>uio_iovdcnt</code>	Specifies the number of already processed <code>iovec</code> elements.
<code>uio_offset</code>	Specifies the character offset on the device performing the I/O.
<code>uio_resid</code>	Specifies the total number of characters remaining in the data area described by the <b>uio</b> structure.
<code>iov_base</code>	Specifies the address of the data area described by the current <code>iovec</code> element.
<code>iov_len</code>	Specifies the length of remaining data area in the buffer described by the current <code>iovec</code> element.

## Execution Environment

The **uiomove** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates successful completion.
ENOMEM	Indicates that there was no room in the buffer.
EIO	Indicates a permanent I/O error file space.
ENOSPC	Indicates insufficient disk space.
EFAULT	Indicates a user location that is not valid.

### Related reference:

“[uphysio Kernel Service](#)” on page 524

“[uio Structure](#)” on page 639

### Related information:

Memory Kernel Services

## unlock\_enable Kernel Service

### Purpose

Unlocks a simple lock if necessary, and restores the interrupt priority.

### Syntax

```
#include <sys/lock_def.h>
```

```
void unlock_enable ( int_pri, lock_addr)
int int_pri;
simple_lock_t lock_addr;
```

## Parameters

Item	Description
<i>int_pri</i>	Specifies the interrupt priority to restore. This must be set to the value returned by the corresponding call to the <b>disable_lock</b> kernel service.
<i>lock_addr</i>	Specifies the address of the lock word to unlock.

## Description

The **unlock\_enable** kernel service unlocks a simple lock if necessary, and restores the interrupt priority, in order to provide optimized thread-interrupt critical section protection for the system on which it is executing. On a multiprocessor system, calling the **unlock\_enable** kernel service is equivalent to calling the **simple\_unlock** and **i\_enable** kernel services. On a uniprocessor system, the call to the **simple\_unlock** service is not necessary, and is omitted. However, you should still pass the valid lock address which was used with the corresponding call to the **disable\_lock** kernel service. Never pass a **NULL** lock address.

## Execution Environment

The **unlock\_enable** kernel service can be called from either the process or interrupt environment.

## Return Values

The **unlock\_enable** kernel service has no return values.

### Related reference:

“disable\_lock Kernel Service” on page 77

“simple\_unlock Kernel Service” on page 476

### Related information:

Understanding Locking

Understanding Interrupts

## unlockl Kernel Service

### Purpose

Unlocks a conventional process lock.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
void unlockl ( lock_word)
lock_t *lock_word;
```

### Parameter

Item	Description
<i>lock_word</i>	Specifies the address of the lock word.

## Description

**Note:** The **unlockl** kernel service is provided for compatibility only and should not be used in new code, which should instead use simple locks or complex locks.

The **unlockl** kernel service unlocks a conventional lock. Only the owner of a lock can unlock it. Once a lock is unlocked, the highest priority thread (if any) which is waiting for the lock is made runnable and may compete again for the lock. If there was at least one process waiting for the lock, the priority of the caller is recomputed. Preempting a System Call discusses how system calls can use locking kernel services when accessing global data.

The **lockl** and **unlockl** services do not maintain a nesting level count. A single call to the **unlockl** service unlocks the lock for the caller. The return code from the **lockl** service should be used to determine when to unlock the lock.

**Note:** The **unlockl** kernel service can be called with interrupts disabled, only if the event or lock word is pinned.

## Execution Environment

The **unlockl** kernel service can be called from the process environment only.

## Return Values

The **unlockl** service has no return values.

## Example

A call to the **unlockl** service can be coded as follows:

```
int lock_ret;          /* return code from lockl() */
extern int lock_word; /* lock word that is external
                      and was initialized to
                      LOCK_AVAIL */
...
/* get lock prior to using resource */
lock_ret = lockl(lock_word, LOCK_SHORT)
/* use resource for which lock was obtained */
...
/* release lock if this was not a nested use */
if ( lock_ret != LOCK_NEST )
    unlockl(lock_word);
```

### Related reference:

“lockl Kernel Service” on page 343

### Related information:

Understanding Locking

## unpin Kernel Service

### Purpose

Unpins the address range in system (kernel) address space.



## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/pin.h>
```

```
int unpin ( addr, length)
caddr addr;
int length;
```

## Parameters

Item	Description
<i>addr</i>	Specifies the address of the first byte to unpin in the system (kernel) address space.
<i>length</i>	Specifies the number of bytes to unpin.

## Description

The **unpin** kernel service decreases the pin count of each page in the address range. When the pin count is 0, the page is not pinned and can be paged out of real memory. Upon finding an unpinned page, the **unpin** service returns the **EINVAL** error code and leaves any remaining pinned pages still pinned.

The **unpin** service can only be called with addresses in the system (kernel) address space. The **xmemunpin** service should be used where the address space might be in either user or kernel space.

## Execution Environment

The **unpin** kernel service can be called from either the process or interrupt environment.

## Return Values

Item	Description
0	Indicates successful completion.
<b>EINVAL</b>	Indicates that the value of the <i>length</i> parameter is negative or 0. Otherwise, the area of memory beginning at the byte specified by the <i>base</i> parameter and extending for the number of bytes specified by the <i>len</i> parameter is not defined. If neither cause is responsible, an unpinned page was specified.

### Related reference:

“pin Kernel Service” on page 408

“xmemunpin Kernel Service” on page 606

### Related information:

Understanding Execution Environments

Memory Kernel Services

## unpincode Kernel Service

### Purpose

Unpins the code and data associated with a loaded object module.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/pin.h>
```

```
int unpincode ( func)
int (*func) ( );
```

## Parameter

Item	Description
<i>func</i>	Specifies an address used to determine the object module to be unpinned. The address is typically that of a function that is exported by this object module.

## Description

The **unpincode** kernel service uses the **ltunpin** kernel service to decrement the pin count for the pages associated with the following items:

- Code associated with the object module
- Data area of the object module that contains the function specified by the *func* parameter

The loader entry for the module is used to determine the size of both the code and the data area.

## Execution Environment

The **unpincode** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates successful completion.
EINVAL	Indicates that the <i>func</i> parameter is not a valid pointer to the function.
EFAULT	Indicates that the calling process does not have access to the area of memory that is associated with the module.

### Related reference:

“unpin Kernel Service” on page 520

### Related information:

Understanding Execution Environments

Memory Kernel Services

## unregister\_HA\_handler Kernel Service

### Purpose

Removes from the kernel the registration of a High Availability Event Handler.

### Syntax

```
#include <sys/high_avail.h>
```

```
int register_HA_handler (ha_handler)  
ha_handler_ext_t * ha_handler;
```

### Parameter

**Item***ha\_handler***Description**

Specifies a pointer to a structure of the type `ha_handler_ext_t` defined in `/usr/include/sys/high_avail.h`. This structure must be identical to the one passed to `register_HA_handler` at the time of registration.

**Description**

The `unregister_HA_handler` kernel service cancels an unconfigured kernel extensions that have registered a high availability event handler, done by the `register_HA_handler` kernel service, so that the kernel extension can be unloaded.

Failure to do so may cause a system crash when a high availability event such as a processor deallocation is initiated due to some hardware fault.

**Execution Environment**

The `unregister_HA_handler` kernel service can be called from the process environment only.

An extension may register the same HAEH  $N$  times ( $N > 1$ ). Although this is considered an incorrect behaviour, no error is reported. The given HAEH will be invoked  $N$  times for each HA event. This handler has to be unregistered as many times as it was registered.

**Return Values****Item**

0

**Description**

Indicates a successful operation.

A non-zero value indicates an error.

**Related reference:**

“register\_HA\_handler Kernel Service” on page 447

**Related information:**

RAS Kernel Services

**untimeout Kernel Service**

**Attention:** This service must not be used because it is not multi-processor safe. The base kernel timer and watchdog services must be used instead. See `talloc` and `w_init` for more information.

**Purpose**

Cancels a pending timer request.

**Syntax**

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
void untimeout ( func, arg)
void (*func)();
caddr_t *arg;
```

**Parameters**

Item	Description
<i>func</i>	Specifies the function that is associated with the timer to be canceled.
<i>arg</i>	Specifies the function argument that is associated with the timer to be canceled.

## Description

The **untimeout** kernel service is not part of the kernel. However, it is a compatibility service that is provided in the **libsys.a** library. To use the **untimeout** service, a kernel extension must be bound with the **libsys.a** library. The **untimeout** service, like the associated kernel libsys services **timeoutcf** and **timeout**, can be bound and used only in the pinned part of a kernel extension or the bottom half of a device driver because these services use interrupt disable for serialization.

The **untimeout** kernel service cancels a specific request that is made with the **timeout** service. The *func* and *arg* parameters must match the parameters that are used in the **timeout** kernel service request that is to be canceled.

Upon return, the specified timer request is canceled, if found. If no timer request matches the *func* and *arg* parameters, no operation is performed.

## Execution Environment

The **untimeout** kernel service can be called from either the process or interrupt environment.

## Return Values

The **untimeout** kernel service has no return values.

### Related reference:

“timeout Kernel Service” on page 490

### Related information:

Timer and Time-of-Day Kernel Services

## uphysio Kernel Service

### Purpose

Performs character I/O for a block device using a **uio** structure.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/buf.h>
#include <sys/uio.h>
```

```
int uphysio (uio, rw, buf_cnt, devno, strat, mincnt, minparms)
struct uio * uio;
int rw;
uint buf_cnt;
dev_t devno;
int (* strat) ( );
int (* mincnt) ( );
void * minparms;
```

### Parameters

Item	Description
<i>uiop</i>	Points to the <b>uiop</b> structure describing the buffer of data to transfer using character-to-block I/O.
<i>rw</i>	Indicates either a read or write operation. A value of <b>B_READ</b> for this flag indicates a read operation. A value of <b>B_WRITE</b> for this flag indicates a write operation.
<i>buf_cnt</i>	Specifies the maximum number of <b>buf</b> structures to use when calling the strategy routine specified by the <i>strat</i> parameter. This parameter is used to indicate the maximum amount of concurrency the device can support and minimize the I/O redrive time. The value of the <i>buf_cnt</i> parameter can range from 1 to 64.
<i>devno</i>	Specifies the major and minor device numbers. With the <b>uphysio</b> service, this parameter specifies the device number to be placed in the <b>buf</b> structure before calling the strategy routine specified by the <i>strat</i> parameter.
<i>strat</i>	Represents the function pointer to the <b>ddstrategy</b> routine for the device.
<i>mincnt</i>	Represents the function pointer to a routine used to reduce the data transfer size specified in the <b>buf</b> structure, as required by the device before the strategy routine is started. The routine can also be used to update extended parameter information in the <b>buf</b> structure before the information is passed to the strategy routine.
<i>minparms</i>	Points to parameters to be used by the <i>mincnt</i> parameter.

## Description

The **uphysio** kernel service performs character I/O for a block device. The **uphysio** service attempts to send to the specified strategy routine the number of **buf** headers specified by the *buf\_cnt* parameter. These **buf** structures are constructed with data from the **uiop** structure specified by the *uiop* parameter.

The **uphysio** service initially transfers data area descriptions from each *iovec* element found in the **uiop** structure into individual **buf** headers. These headers are later sent to the strategy routine. The **uphysio** kernel service tries to process as many data areas as the number of **buf** headers permits. It then invokes the strategy routine with the list of **buf** headers.

### Preparing Individual **buf** Headers

The routine specified by the *mincnt* parameter is called before the **buf** header, built from an *iovec* element, is added to the list of **buf** headers to be sent to the strategy routine. The *mincnt* parameter is passed a pointer to the **buf** header along with the *minparms* pointer. This arrangement allows the *mincnt* parameter to tailor the length of the data transfer described by the **buf** header as required by the device performing the I/O. The *mincnt* parameter can also optionally modify certain device-dependent fields in the **buf** header.

When the *mincnt* parameter returns with no error, an attempt is made to pin the data buffer described by the **buf** header. If the pin operation fails due to insufficient memory, the data area described by the **buf** header is reduced by half. The **buf** header is again passed to the *mincnt* parameter for modification before trying to pin the reduced data area.

This process of downsizing the transfer specified by the **buf** header is repeated until one of the three following conditions occurs:

- The pin operation succeeds.
- The *mincnt* parameter indicates an error.
- The data area size is reduced to 0.

When insufficient memory indicates a failed pin operation, the number of **buf** headers used for the remainder of the operation is reduced to 1. This is because trying to pin multiple data areas simultaneously under these conditions is not desirable.

If the user has not already obtained cross-memory descriptors, further processing is required. (The `uio_segflg` field in the `uio` structure indicates whether the user has already initialized the cross-memory descriptors. The `usr/include/sys/uio.h` file contains information on possible values for this flag.)

When the data area described by the `buf` header has been successfully pinned, the `uphysio` service verifies user access authority for the data area. It also obtains a cross-memory descriptor to allow the device driver interrupt handler limited access to the data area.

### Calling the Strategy Routine

After the `uphysio` kernel service obtains a cross-memory descriptor to allow the device driver interrupt handler limited access to the data area, the `buf` header is then put on a list of `buf` headers to be sent to the strategy routine specified by the `strat` parameter.

The strategy routine specified by the `strat` parameter is called with the list of `buf` headers when:

- The list reaches the number of `buf` structures specified by the `buf_cnt` parameter.
- The data area described by the `uio` structure has been completely described by `buf` headers.

The `buf` headers in the list are chained together using the `av_back` and `av_forw` fields before they are sent to the strategy routine.

### Waiting for `buf` Header Completion

When all available `buf` headers have been sent to the strategy routine, the `uphysio` service waits for one or more of the `buf` headers to be marked complete. The `IODONE` handler is used to wake up the `uphysio` service when it is waiting for completed `buf` headers from the strategy routine.

When the `uphysio` service is notified of a completed `buf` header, the associated data buffer is unpinned and the cross-memory descriptor is freed. (However, the cross-memory descriptor is freed only if the user had not already obtained it.) An error is detected on the data transfer under the following conditions:

- The completed `buf` header has a nonzero `b_resid` field.
- The `b_flags` field has the `B_ERROR` flag set.

When an error is detected by the `uphysio` service, no new `buf` headers are sent to the strategy routine.

The `uphysio` service waits for any `buf` headers already sent to the strategy routine to be completed and then returns an error code to the caller. If no errors are detected, the `buf` header and any other completed `buf` headers are again used to send more data transfer requests to the strategy routine as they become available. This process continues until all data described in the `uio` structure has been transferred or until an error has been detected.

The `uphysio` service returns to the caller when:

- All `buf` headers have been marked complete by the strategy routine.
- All data specified by the `uio` structure has been transferred.

The `uphysio` service also returns an error code to the caller if an error is detected.

### Error Detection by the `uphysio` Kernel Service

When it detects an error, the `uphysio` kernel service reports the error that was detected closest to the start of the data area described by the `uio` structure. No additional `buf` headers are sent to the strategy routine. The `uphysio` kernel service waits for all `buf` headers sent to the strategy routine to be marked complete.

However, additional `buf` headers may have been sent to the strategy routine between these two events:

- After the strategy routine detects the error.
- Before the **uphysio** service is notified of the error condition in the completed **buf** header.

When errors occur, various fields in the returned **uio** structure may or may not reflect the error. The **uio\_iov** and **uio\_iovcnt** fields are not updated and contain their original values.

The **uio\_resid** and **uio\_offset** fields in the returned **uio** structure indicate the number of bytes transferred by the strategy routine according to the sum of all (the **b\_bcount** field minus the **b\_resid** fields) fields in the **buf** headers processed by the strategy routine. These headers include the **buf** header indicating the error nearest the start of the data area described by the original **uio** structure. Any data counts in **buf** headers completed after the detection of the error are not reflected in the returned **uio** structure.

## Execution Environment

The **uphysio** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates successful completion.
ENOMEM	Indicates that no memory is available for the required <b>buf</b> headers.
EAGAIN	Indicates that the operation fails due to a temporary insufficient resource condition.
EFAULT	Indicates that the <b>uio_segflg</b> field indicated user space and that the user does not have authority to access the buffer.
EIO or the <b>b_error</b> field in a <b>buf</b> header	Indicates an I/O error in a <b>buf</b> header processed by the strategy routine.
Return code from the <b>mincnt</b> parameter	Indicates that the return code from the <b>mincnt</b> parameter if the routine returned with a nonzero return code.

### Related reference:

“uphysio Kernel Service **mincnt** Routine”

“**buf** Structure” on page 615

“**uio** Structure” on page 639

## uphysio Kernel Service **mincnt** Routine Purpose

Tailors a **buf** data transfer request to device-dependent requirements.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/buf.h>

int mincnt ( bp, minparms)
struct buf *bp;
void *minparms;
```

### Parameters

Item	Description
<i>bp</i>	Points to the <b>buf</b> structure to be tailored.
<i>minparms</i>	Points to parameters.

## Description

Only the following fields in the **buf** header sent to the routine specified by the **uphysio** kernel service *mincnt* parameter can be modified by that routine:

- *b\_bcount*
- *b\_work*
- *b\_options*

The *mincnt* parameter cannot modify any other fields without the risk of error. If the *mincnt* parameter determines that the **buf** header cannot be supported by the target device, the routine should return a nonzero return code. This stops the **buf** header and any additional **buf** headers from being sent to the **ddstrategy** routine.

The **uphysio** kernel service waits for all **buf** headers already sent to the strategy routine to complete and then returns with the return code from the *mincnt* parameter.

### Related reference:

“uphysio Kernel Service” on page 524

## uprintf Kernel Service

### Purpose

Submits a request to print a message to the controlling terminal of a process.

### Syntax

```
#include <sys/uprintf.h>
```

```
int uprintf ( Format [,
             Value, ...])
char *Format;
```



## Parameters

Item	Description
<i>Format</i>	<p>Specifies a character string containing either or both of two types of objects:</p> <ul style="list-style-type: none"><li>• Plain characters, which are copied to the message output stream.</li><li>• Conversion specifications, each of which causes 0 or more items to be retrieved from the <i>Value</i> parameter list. Each conversion specification consists of a % (percent sign) followed by a character that indicates the type of conversion to be applied:</li></ul> <p>%        Performs no conversion. Prints %.</p> <p>d, i      Accepts an integer <i>Value</i> and converts it to signed decimal notation.</p> <p>u        Accepts an integer <i>Value</i> and converts it to unsigned decimal notation.</p> <p>o        Accepts an integer <i>Value</i> and converts it to unsigned octal notation.</p> <p>x        Accepts an integer <i>Value</i> and converts it to unsigned hexadecimal notation.</p> <p>s        Accepts a <i>Value</i> as a string (character pointer), and characters from the string are printed until a \0 (null character) is encountered. <i>Value</i> must be non-null and the maximum length of the string is limited to <b>UP_MAXSTR</b> characters.</p> <p>Field width or precision conversion specifications are not supported.</p> <p>The following constants are defined in the <code>/usr/include/sys/uprintf.h</code> file:</p> <ul style="list-style-type: none"><li>– <b>UP_MAXSTR</b></li><li>– <b>UP_MAXARGS</b></li><li>– <b>UP_MAXCAT</b></li><li>– <b>UP_MAXMSG</b></li></ul> <p>The <i>Format</i> string may contain from 0 to the number of conversion specifications specified by the <b>UP_MAXARGS</b> constant. The maximum length of the <i>Format</i> string is the number of characters specified by the <b>UP_MAXSTR</b> constant. <i>Format</i> must be non-null.</p> <p>The maximum length of the constructed kernel message is limited to the number of characters specified by the <b>UP_MAXMSG</b> constant. Messages larger than the number of characters specified by the <b>UP_MAXMSG</b> constant are discarded.</p>
<i>Value</i>	<p>Specifies, as an array, the value to be converted. The number, type, and order of items in the <i>Value</i> parameter list should match the conversion specifications within the <i>Format</i> string.</p>

## Description

The **uprintf** kernel service submits a kernel message request. Once the request has been successfully submitted, the **uprintfd** daemon constructs the message based on the *Format* and *Value* parameters of the request. The **uprintfd** daemon then writes the message to the process' controlling terminal.

## Execution Environment

The **uprintf** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
ENOMEM	Indicates that memory is not available to buffer the request.
ENODEV	Indicates that a controlling terminal does not exist for the process.
ESRCH	Indicates that the <b>uprintfd</b> daemon is not active. No requests may be submitted.
EINVAL	Indicates that a string <i>Value</i> string pointer is null or the string <i>Value</i> parameter is greater than the number of characters specified by the <b>UP_MAXSTR</b> constant.
EINVAL	Indicates one of the following: <ul style="list-style-type: none"><li>• <i>Format</i> string pointer is null.</li><li>• Number of characters in the <i>Format</i> string is greater than the number specified by the <b>UP_MAXSTR</b> constant.</li><li>• Number of conversion specifications contained within the <i>Format</i> string is greater than the number specified by the <b>UP_MAXARGS</b> constant.</li></ul>

**Related reference:**

“NLuprintf Kernel Service” on page 386

**Related information:**

uprintfd command

Process and Exception Management Kernel Services

**ureadc Kernel Service****Purpose**

Writes a character to a buffer described by a **uio** structure.

**Syntax**

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/uio.h>
```

```
int ureadc ( c, uio)
int c;
struct uio *uio;
```

**Parameters**

Item	Description
<i>c</i>	Specifies a character to be written to the buffer.
<i>uio</i>	Points to a <b>uio</b> structure describing the buffer in which to place a character.

**Description**

The **ureadc** kernel service writes a character to a buffer described by a **uio** structure. Device driver top half routines, especially character device drivers, frequently use the **ureadc** kernel service to transfer data into a user area.

The *uio\_resid* and *uio\_iovcnt* fields in the **uio** structure describing the data area must be greater than 0. If these fields are not greater than 0, an error is returned. The *uio\_segflg* field in the **uio** structure is used to indicate whether the data is being written to a user- or kernel-data area. It is also used to indicate if the caller requires cross-memory operations and has provided the required cross-memory descriptors. The values for the flag are defined in the `/usr/include/sys/uio.h` file.

If the data is successfully written, the following fields in the **uio** structure are updated:

Field	Description
<i>uio_iov</i>	Specifies the address of current <i>iovec</i> element to use.
<i>uio_xmem</i>	Specifies the address of current <i>xmem</i> element to use (used for cross-memory copy).
<i>uio_iovcnt</i>	Specifies the number of remaining <i>iovec</i> elements.
<i>uio_iovdcnt</i>	Specifies the number of <i>iovec</i> elements already processed.
<i>uio_offset</i>	Specifies the character offset on the device from which data is read.
<i>uio_resid</i>	Specifies the total number of characters remaining in the data area described by the <i>uio</i> structure.
<i>iov_base</i>	Specifies the address of the next available character in the data area described by the current <i>iovec</i> element.
<i>iov_len</i>	Specifies the length of remaining data area in the buffer described by the current <i>iovec</i> element.

**Execution Environment**

The **ureadc** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates successful completion.
ENOMEM	Indicates that there is no room in the buffer.
EFAULT	Indicates that the user location is not valid for one of these reasons: <ul style="list-style-type: none"><li>• The <code>uio_segflg</code> field indicates user space and the base address (<code>iov_base</code> field) points to a location outside of the user address space.</li><li>• The user does not have sufficient authority to access the location.</li><li>• An I/O error occurs while accessing the location.</li></ul>

### Related reference:

“uio move Kernel Service” on page 517

“uwritec Kernel Service”

### Related information:

Memory Kernel Services

## uwritec Kernel Service

### Purpose

Retrieves a character from a buffer described by a **uio** structure.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/uio.h>
```

```
int uwritec ( uiop)
struct uio *uiop;
```

### Parameter

Item	Description
<i>uiop</i>	Points to a <b>uio</b> structure describing the buffer from which to read a character.

### Description

The **uwritec** kernel service reads a character from a buffer described by a **uio** structure. Device driver top half routines, especially character device drivers, frequently use the **uwritec** kernel service to transfer data out of a user area. The `uio_resid` and `uio_iovcnt` fields in the **uio** structure must be greater than 0 or an error is returned.

The `uio_segflg` field in the **uio** structure indicates whether the data is being read out of a user- or kernel-data area. This field also indicates whether the caller requires cross-memory operations and has provided the required cross-memory descriptors. The values for this flag are defined in the `/usr/include/sys/uio.h` file.

If the data is successfully read, the following fields in the **uio** structure are updated:

Field	Description
<code>uio_iov</code>	Specifies the address of the current <code>iovec</code> element to use.
<code>uio_xmem</code>	Specifies the address of the current <code>xmem</code> element to use (used for cross-memory copy).
<code>uio_iovcnt</code>	Specifies the number of remaining <code>iovec</code> elements.
<code>uio_iovdcnt</code>	Specifies the number of <code>iovec</code> elements already processed.
<code>uio_offset</code>	Specifies the character offset on the device to which data is written.
<code>uio_resid</code>	Specifies the total number of characters remaining in the data area described by the <code>uio</code> structure.
<code>iov_base</code>	Specifies the address of the next available character in the data area described by the current <code>iovec</code> element.
<code>iov_len</code>	Specifies the length of the remaining data in the buffer described by the current <code>iovec</code> element.

## Execution Environment

The `uwritec` kernel service can be called from the process environment only.

## Return Values

Upon successful completion, the `uwritec` service returns the character it was sent to retrieve.

Item	Description
-1	Indicates that the buffer is empty or the user location is not valid for one of these three reasons: <ul style="list-style-type: none"> <li>The <code>uio_segflg</code> field indicates user space and the base address (<code>iov_base</code> field) points to a location outside of the user address space.</li> <li>The user does not have sufficient authority to access the location.</li> <li>An I/O error occurred while the location was being accessed.</li> </ul>

### Related reference:

“`uio`move Kernel Service” on page 517

“`u`physio Kernel Service” on page 524

“`u`readc Kernel Service” on page 530

## V

The following kernel services begin with the letter v.

### validate\_pag or validate\_pag64 Kernel Service

#### Purpose

Validates the Process Authentication Group (PAG) value.

#### Syntax

```
#include <sys/cred.h>
```

```
int validate_pag ( type, pg, npags )
int type;
struct paglist pg[];
int npags;
```

```
int validate_pag64 ( type, pg, npags )
int type;
struct paglist64 pg[];
int npags;
```

#### Parameters

Item	Description
<i>type</i>	PAG type to validate
<i>pg</i>	PAG list (must be in pinned memory)
<i>npags</i>	Number of PAGs to validate

## Description

The **validate\_pag** or **validate\_pag64** kernel service validates the PAGs specified in *pg*. These services support the garbage collection of data structures by kernel extensions associated with PAGs. These structures are associated with a **set\_pag** interface process. PAG values are inherited from parent to child across the **fork** system call, so one kernel extension structure can map to many processes. This routine is required to synchronize the execution of forks so that the process table can be scanned to identify a particular PAG. The **validate\_pag** and **validate\_pag64** kernel services cannot be used simultaneously with the **set\_pag** interface. The application is required to provide this synchronization.

The value of *type* must be a defined PAG ID. The PAG ID for the Distributed Computing Environment (DCE) is 0. The *pg* parameter must be a valid, referenced PAG list in pinned memory.

## Execution Environment

The **validate\_pag** and **validate\_pag64** kernel services can be called from the process environment only.

## Return Values

A value of 0 is returned upon successful completion. Upon failure, a -1 is returned and **errno** is set to a value that explains the error.

## Error Codes

The **validate\_pag** and **validate\_pag64** kernel services fail if the following condition is true:

Item	Description
EINVAL	Invalid PAG specification

## Related Information

Security Kernel Services in *Kernel Extensions and Device Support Programming Concepts*.

### Related information:

Security Kernel Services

## vec\_clear Kernel Service

### Purpose

Removes a virtual interrupt handler.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
void vec_clear ( levsublev)
int levsublev;
```

## Parameter

Item	Description
<i>levsublev</i>	Represents the value returned by <b>vec_init</b> kernel service when the virtual interrupt handler was defined.

## Description

The **vec\_clear** kernel service is not part of the base kernel but is provided by the device queue management kernel extension. This queue management kernel extension must be loaded into the kernel before loading any kernel extensions referencing these services.

The **vec\_clear** kernel service removes the association between a virtual interrupt handler and the virtual interrupt level and sublevel that was assigned by the **vec\_init** kernel service. The virtual interrupt handler at the sublevel specified by the *levsublev* parameter no longer registers upon return from this routine.

## Execution Environment

The **vec\_clear** kernel service can be called from the process environment only.

## Return Values

The **vec\_clear** kernel service has no return values. If no virtual interrupt handler is registered at the specified sublevel, no operation is performed.

### Related reference:

“**vec\_init** Kernel Service”

## **vec\_init** Kernel Service

### Purpose

Defines a virtual interrupt handler.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
int vec_init ( level, routine, arg)
int level;
void (*routine) ();
int arg;
```

### Parameters

Item	Description
<i>level</i>	Specifies the virtual interrupt level. This level value is not used by the <b>vec_init</b> kernel service and implies no relative priority. However, it is returned with the sublevel assigned for the registered virtual interrupt handler.
<i>routine</i>	Identifies the routine to call when a virtual interrupt occurs on a given interrupt sublevel.
<i>arg</i>	Specifies a value that is passed to the virtual interrupt handler.

## Description

The **vec\_init** kernel service is not part of the base kernel but provided by the device queue management kernel extension. This queue management kernel extension must be loaded into the kernel before loading any kernel extensions referencing these services.

The **vec\_init** kernel service associates a virtual interrupt handler with a level and sublevel. This service searches the available sublevels to find the first unused one. The *routine* and *arg* parameters are used to initialize the open sublevel. The **vec\_init** kernel service then returns the level and assigned sublevel.

There is a maximum number of available sublevels. If this number is exceeded, the **vec\_init** service halts the system. This service should be called to initialize a virtual interrupt before any device queues using the virtual interrupt are created.

The *level* parameter is not used by the **vec\_init** service. It is provided for compatibility reasons only. However, its value is passed back intact with the sublevel.

## Execution Environment

The **vec\_init** kernel service can be called from the process environment only.

## Return Values

The **vec\_init** kernel service returns a value that identifies the virtual interrupt level and assigned sublevel. The low-order 8 bits of this value specify the sublevel, and the high-order 8 bits specify the level. The **attachq** kernel service uses the same format. This level value is the same value as that supplied by the *level* parameter.

## vfsrele Kernel Service

### Purpose

Releases all resources associated with a virtual file system.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
int vfsrele ( vfsp)
struct vfs *vfsp;
```

### Parameter

Item	Description
<i>vfsp</i>	Points to a virtual file system structure.

### Description

The **vfsrele** kernel service releases all resources associated with a virtual file system.

When a file system is unmounted, the **VFS\_UNMOUNTED** flag is set in the **vfs** structure, indicating that it is no longer valid to do path name-related operations within the file system. When this flag is set and a **vnop\_rele** v-node operation releases the last active v-node within the file system, the **vnop\_rele** v-node implementation must call the **vfsrele** kernel service to complete the deallocation of the **vfs** structure.

## Execution Environment

The **vfsrele** kernel service can be called from the process environment only.

## Return Values

The `vfsrele` kernel service always returns a value of 0.

### Related information:

Virtual File System Overview

Virtual File System (VFS) Kernel Services

## vm\_att Kernel Service

### Purpose

Maps a specified virtual memory object to a region in the current address space.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
caddr_t vm_att ( vmhandle, offset)
vmhandle_t vmhandle;
caddr_t offset;
```

### Parameters

Item	Description
<i>vmhandle</i>	Specifies the handle for the virtual memory object to be mapped.
<i>offset</i>	Specifies the offset in the virtual memory object and region.

### Description

The `vm_att` kernel service performs the following tasks:

- Selects an unallocated region in the current address space and allocates it.
- Maps the virtual memory object specified by the *vmhandle* parameter with the access permission specified in the handle.
- Constructs the address in the current address space corresponding to the offset in the virtual memory object and region.

The `vm_att` kernel service assumes an address space model of fixed-size virtual memory objects and address space regions.

**Attention:** If there are no more free regions, this call cannot complete and calls the `panic` kernel service.

### Execution Environment

The `vm_att` kernel service can be called from either the process or interrupt environment.

### Return Values

The `vm_att` kernel service returns the address that corresponds to the *offset* parameter in the address space.

### Related reference:

“`vm_det` Kernel Service” on page 537

### Related information:



## **vm\_cflush Kernel Service**

### **Purpose**

Flushes the processor's cache for a specified address range.

### **Syntax**

```
#include <sys/types.h>  
#include <sys/errno.h>  
#include <sys/vmuser.h>
```

```
void vm_cflush ( eaddr, nbytes)  
caddr_t eaddr;  
int nbytes;
```

### **Parameters**

Item	Description
<i>eaddr</i>	Specifies the starting address of the specified range.
<i>nbytes</i>	Specifies the number of bytes in the address range. If this parameter is negative or 0, no lines are invalidated.

### **Description**

The **vm\_cflush** kernel service writes to memory all modified cache lines that intersect the address range (*eaddr*, *eaddr* + *nbytes* - 1). The *eaddr* parameter can have any alignment in a page.

The **vm\_cflush** kernel service can only be called with addresses in the system (kernel) address space.

### **Execution Environment**

The **vm\_cflush** kernel service can be called from both the interrupt and the process environment.

### **Return Values**

The **vm\_cflush** kernel service has no return values.

#### **Related information:**

Memory Kernel Services

Understanding Virtual Memory Manager Interfaces

## **vm\_det Kernel Service**

### **Purpose**

Unmaps and deallocates the region in the current address space that contains a given address.

### **Syntax**

```
#include <sys/types.h>  
#include <sys/errno.h>  
#include <sys/vmuser.h>
```

```
void vm_det ( eaddr)  
caddr_t eaddr;
```

## Parameter

Item	Description
<i>eaddr</i>	Specifies the effective address in the current address space. The region containing this address is to be unmapped and deallocated.

## Description

The **vm\_det** kernel service unmaps the region containing the *eaddr* parameter and deallocates the region, adding it to the free list for the current address space.

The **vm\_det** kernel service assumes an address space model of fixed-size virtual memory objects and address space regions.

**Attention:** If the region is not mapped, or a system region is referenced, the system will halt.

## Execution Environment

The **vm\_det** kernel service can be called from either the process or interrupt environment.

### Related reference:

“vm\_att Kernel Service” on page 536

### Related information:

Memory Kernel Services

Understanding Virtual Memory Manager Interfaces

## vm\_flushp Kernel Service

### Purpose

Flushes the specified range of pages.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
int vm_flushp ( sid, pfirst, npages)
vmid_t sid;
vpn_t pfirst;
vpn_t npages;
```

### Parameters

Item	Description
<i>sid</i>	Identifies the base segment.
<i>pfirst</i>	The first page number within the range.
<i>npages</i>	The number of pages to flush starting from the <i>pfirst</i> value. All pages must be in the same segment.

## Description

The **vm\_flushp** kernel service routine initiates page-out for the specified page range in the virtual memory object. I/O is initiated for the modified pages only. If page-out is initiated, or the pages are currently undergoing page I/O, then they are flagged to have their page frames released upon completion. If the pages are not modified, their page frames are immediately released.

The caller can wait for the completion of I/O initiated by this and prior calls by calling the `vms_iowait` kernel service.

**Note:** The `vm_flushp` subroutine is not supported for use on large pages.

## Execution Environment

The `vm_flushp` kernel service can be called from the process environment only.

This is intended for files, and might not be called for working storage segments.

## Return Values

Item	Description
0	Indicates the completion of the flush operation.
EINVAL	Indicates one of the following errors: <ul style="list-style-type: none"><li>• <code>pfirst = 0</code> and <code>npages = 0</code>.</li><li>• <code>pfirst &lt; 0</code>.</li><li>• <code>npages &lt; 0</code>.</li><li>• Page interval not in one segment.</li><li>• Invalid <code>sid</code> parameter.</li><li>• Invalid segment type.</li></ul>

### Related reference:

“`vm_wripte` Kernel Service” on page 579

“`vm_invalidatep` Kernel Service” on page 544

### Related information:

Understanding Virtual Memory Manager Interfaces

## vm\_galloc Kernel Service

### Purpose

Allocates a region of global memory in the 64-bit kernel.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
int vm_galloc (int type, vm_size_t size, ulong * eaddr)
```

### Description

The `vm_galloc` kernel service allocates memory from the kernel global memory pool on the 64-bit kernel. The allocation size is rounded up to the nearest 4K boundary. The default page protection key for global memory segments is 00 unless overridden with the `V_UREAD` flag.

The `type` field may have the following values, which may be combined:

Item	Description
V_WORKING	Required. Creates a working storage segment.
V_SYSTEM	The new allocation is a global system area that does not belong to any application. Storage reference errors to this area will result in system crashes.
V_UREAD	Overrides the default page protection of 00 and creates the new region with a default page protection of 01.
V_NOEXEC	Pages in the region will have no-execute protection by default. Only supported on POWER4 and later hardware.

The **vm\_galloc** kernel service is intended for subsystems that have large data structures for which **xmalloc** is not the best choice for management. The kernel **xmalloc** heap itself does reside in global memory.

## Parameters

Item	Description
<i>type</i>	Flags that may be specified to control the allocation.
<i>size</i>	Specifies the size, in bytes, of the desired allocation.
<i>eaddr</i>	Pointer to where <b>vm_galloc</b> will return the start address of the allocated storage.

## Execution Environment

The **vm\_galloc** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Successful completion. A new region was allocated, and its start address is returned at the address specified by the <b>eaddr</b> parameter.
EINVAL	Invalid size or type specified.
ENOSPC	Not enough space in the <b>galloc</b> heap to perform the allocation.
ENOMEM	Insufficient resources available to satisfy the request.

### Related reference:

“vm\_gfree Kernel Service”

### Related information:

Memory Kernel Services

Understanding Virtual Memory Manager Interfaces

## vm\_gfree Kernel Service

### Purpose

Frees a region of global memory in the kernel previously allocated with the **vm\_galloc** kernel service.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
int vm_gfree (ulong eaddr, vm_size_t size)
```

## Description

The `vm_gfree` kernel service frees up a global memory region previously allocated with the `vm_galloc` kernel service. The start address and size must exactly match what was previously allocated by the `vm_galloc` kernel service. It is not valid to free part of a previously allocated region in the `vm_galloc` area.

Any I/O to or from the region being freed up must be quiesced before calling the `vm_gfree` kernel service.

## Parameters

Item	Description
<i>eaddr</i>	Start address of the region to free.
<i>size</i>	Size in bytes of the region to free.

## Execution Environment

The `vm_gfree` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Successful completion. The region was freed.
EINVAL	Invalid size or start address specified. This could mean that the region is out of range of the <code>vm_galloc</code> heap, was not previously allocated with <code>vm_galloc</code> , or does not exactly match a previous allocation from <code>vm_galloc</code> .

### Related reference:

“`vm_galloc` Kernel Service” on page 539

### Related information:

Memory Kernel Services

Understanding Virtual Memory Manager Interfaces

## vm\_guatt Kernel Service

### Purpose

Attaches an area of global kernel memory to the current process's address space.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
int vm_guatt (kaddr, size, key, flags, uaddr)
void * kaddr;
vm_size_t size;
vmkey_t key;
long flags;
void ** uaddr;
```

### Parameters

Item	Description
<i>kaddr</i>	Kernel address to be attached (returned from <code>vm_galloc</code> when the global memory was allocated).
<i>size</i>	Length of the region to be inserted into the process address space, in bytes.
<i>key</i>	Protection key that the process will use when accessing the attached region.
<i>flags</i>	Type of <code>vm_guatt</code> operation; must be set to <code>VU_ANYWHERE</code> .
<i>uaddr</i>	Pointer to user space address where the region was attached by <code>vm_guatt</code> . The location pointed to by <i>uaddr</i> ( <i>*uaddr</i> ) must be null when the <code>vm_guatt</code> call is made.

## Description

`vm_guatt` is a kernel service used to attach a region of global kernel memory that was allocated with `vm_galloc` to a process's address space. If the call is successful, the address in the process address space where the memory was attached is returned in the location pointed to by *uaddr*.

*key* can be set to `VM_PRIV` or `VM_UNPRIV`. If it is set to `VM_PRIV`, the process will be able to read and write the attached region. If it is set to `VM_UNPRIV`, the process will not be able to write the region and will only be able to read it if the `vm_galloc` of the region was done with the `V_UREAD` flag on.

`vm_guatt` attachments are not inherited across a process fork.

## Execution Environment

The `vm_guatt` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
EINVAL	Indicates one of the following errors: <ul style="list-style-type: none"> <li><i>flags</i> or <i>key</i> is not set to a valid value, <i>size</i> is 0, or the value pointed to by <i>uaddr</i> is non-NULL.</li> <li>Region indicated by <i>kaddr</i> and <i>size</i> does not lie within a region previously allocated by <code>vm_galloc</code>.</li> </ul>

## Implementation Specifics

The `vm_guatt` kernel service is part of Base Operating System (BOS) Runtime.

### Related reference:

“`vm_galloc` Kernel Service” on page 539

“`vm_gudet` Kernel Service”

### Related information:

Memory Kernel Services

## `vm_gudet` Kernel Service

### Purpose

Removes a region attached with `vm_guatt` from the current process's address space.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
int vm_gudet (kaddr, uaddr, size, flags)
```

```
void * kaddr;  
void * uaddr;  
vm_size_t size;  
long flags;
```

## Parameters

Item	Description
<i>kaddr</i>	Kernel address attached by <code>vm_guatt</code> .
<i>uaddr</i>	Location in the process address space where the kernel region was attached.
<i>size</i>	Length of the attached region, in bytes.
<i>flags</i>	Type of <code>vm_gudet</code> operation, must be <code>VU_ANYWHERE</code> .

## Description

`vm_gudet` is a kernel service that detaches a region of global kernel memory that was attached by `vm_guatt`. This memory must still be allocated, detaching a region after it has been deallocated with `vm_gfree` is an error. If the detach is successful, the global kernel memory region at *kaddr* will no longer be addressable at *uaddr* by the calling process.

## Execution Environment

The `vm_gudet` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	User address detached successfully.
EINVAL	Indicates one of the following errors: <ul style="list-style-type: none"><li>Invalid flags.</li><li>Region indicated by <i>kaddr</i> and <i>size</i> does not lie within a region allocated by <code>vm_galloc</code>.</li></ul>

## Implementation Specifics

The `vm_gudet` kernel service is part of Base Operating System (BOS) Runtime.

### Related reference:

“`vm_galloc` Kernel Service” on page 539

“`vm_gfree` Kernel Service” on page 540

“`vm_guatt` Kernel Service” on page 541

### Related information:

Memory Kernel Services

## `vm_handle` Kernel Service

### Purpose

Constructs a virtual memory handle for mapping a virtual memory object with a specified access level.

### Syntax

```
#include <sys/types.h>  
#include <sys/errno.h>  
#include <sys/vmuser.h>
```

```
vmhandle_t vm_handle ( vmid, key )  
vmid_t vmid;  
int key;
```

## Parameters

Item	Description
<i>vmid</i>	Specifies a virtual memory object identifier, as returned by the <b>vms_create</b> kernel service.
<i>key</i>	Specifies an access key. This parameter has a 1 value for limited access and a 0 value for unlimited access, respectively.

## Description

The **vm\_handle** kernel service constructs a virtual memory handle for use by the **vm\_att** kernel service. The handle identifies the virtual memory object specified by the *vmid* parameter and contains the access key specified by the *key* parameter.

A virtual memory handle is used with the **vm\_att** kernel service to map a virtual memory object into the current address space.

The **vm\_handle** kernel service assumes an address space model of fixed-size virtual memory objects and address space regions.

## Execution Environment

The **vm\_handle** kernel service can be called from the process environment only.

## Return Values

The **vm\_handle** kernel service returns a virtual memory handle type.

### Related reference:

“vm\_att Kernel Service” on page 536

“vms\_create Kernel Service” on page 572

### Related information:

Understanding Virtual Memory Manager Interfaces

## vm\_invalidatep Kernel Service

### Purpose

Releases page frames in the specified range for a non-journaled persistent segment or client segment.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
int vm_invalidatep ( sid, pfirst, npages)
```

```
vmid_t sid;
```

```
vpn_t pfirst;
```

```
ulong npages;
```

### Parameters



Item	Description
<i>sid</i>	Identifies the base segment.
<i>pfirst</i>	The first page number within the range.
<i>npages</i>	The number of pages to invalidate starting from the <i>pfirst</i> value. All pages must be in the same segment.

## Description

The **vm\_invalidatep** kernel service routine discards any page frames associated with the virtual memory object in the specified page range.

If a page within the specified range is found in page-in or page-out state, then the thread is synchronously put to sleep until the page I/O completes. When the I/O is complete, any memory-resident page frame is then freed.

**Note:** The **vm\_invalidatep** subroutine is not supported for use on large pages.

## Execution Environment

The **vm\_invalidatep** kernel service can be called from the process environment only.

This is intended for files, and might not be called for working storage segments.

## Return Values

Item	Description
0	Indicates the completion of the invalidating operations.
EINVAL	Indicates one of the following errors: <ul style="list-style-type: none"> <li>• <i>pfirst</i> &lt; 0.</li> <li>• <i>npages</i> &lt; 0.</li> <li>• Page interval not in one segment.</li> <li>• Invalid <i>sid</i> parameter.</li> <li>• Invalid segment type.</li> </ul>

### Related reference:

“vm\_wriep Kernel Service” on page 579

“vms\_iowait, vms\_iowaitf Kernel Services” on page 574

### Related information:

Understanding Virtual Memory Manager Interfaces

## vm\_ioaccessp Kernel Service

### Purpose

Initiates asynchronous page-in or page-out for the range of pages specified.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
int vm_ioaccessp ( bsid, pfirst, npages, modifier)
vmid_t bsid;
vpn_t pfirst;
vpn_t npages;
uint modifier;
```

## Parameters

Item	Description
<i>bsid</i>	Identifies the base segment.
<i>first</i>	The first page number within the range.
<i>npages</i>	The number of pages to access starting from the <i>first</i> value. All pages must be in the same segment.
<i>modifier</i>	Flags passed in by the user. These flags are detailed below.

## Description

The **vm\_ioaccessp** kernel service routine enables a client file system with a thread-level strategy routine to access the pages specified. This call is strictly advisory and might return without having done anything. If you want to actually move the data, call the **vm\_uiomove** kernel service. If you want to pre-page the target, then call the **vm\_readp** kernel service.

The flags passed in through the *modifier* parameter determine what type of action taken by the **vm\_ioaccessp** kernel service. For details of each flag's purpose, see the table below.

The flags carry certain restrictions. You cannot request both a make and a flush operation. Also, if the **VM\_IOACCESSP\_WAITONLY** flag is declared then you must specify at least one type of wait operation. Finally, you cannot request a make or a flush operation if the **VM\_IOACCESSP\_WAITONLY** flag is declared.

## Flags

Value	Name	Purpose
0x0001	<b>VM_IOACCESSP_MAKE</b>	Creates new pages in the page-in state in the specified range. Can only make up to 1MB of pages.
0x0002	<b>VM_IOACCESSP_FLUSH</b>	Flushes pages in the specified range.
0x0004	<b>VM_IOACCESSP_PGINWAIT</b>	If a page in the specified range is in page-in state, then block until page-in is complete.
0x0008	<b>VM_IOACCESSP_PGOUTWAIT</b>	If a page in the specified range is in page-out state, then block until page-out is complete.
0x0010	<b>VM_IOACCESSP_WAITONLY</b>	Returns once the specified wait is complete. The <b>VM_IOACCESSP_PGINWAIT</b> flag and the <b>VM_IOACCESSP_PGOUTWAIT</b> flag must also be specified.

## Execution Environment

The **vm\_ioaccessp** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates the completion of the I/O access operations.

Item	Description
EINVAL	<p>Indicates one of the following errors:</p> <ul style="list-style-type: none"> <li>• <i>pfirst</i> = 0 and <i>npages</i> = 0.</li> <li>• <i>pfirst</i> &lt; 0.</li> <li>• <i>npages</i> &lt; 0.</li> <li>• Page interval not in one segment.</li> <li>• Invalid <i>sid</i> parameter.</li> <li>• Page make requests &gt; 1 MB.</li> <li>• Not a client file system.</li> <li>• Unsupported flag used.</li> <li>• Both the <b>VM_IOACCESSP_MAKE</b> and the <b>VM_IOACCESSP_FLUSH</b> flags are set.</li> <li>• The <b>VM_IOACCESSP_WAITONLY</b> flag is set and the <b>VM_IOACCESSP_PGINWAIT</b> flag or the <b>VM_IOACCESSP_PGOUTWAIT</b> flag is not set.</li> <li>• The <b>VM_IOACCESSP_WAITONLY</b> flag and the <b>VM_IOACCESSP_MAKE</b> flag or the <b>VM_IOACCESSP_FLUSH</b> flag are set.</li> </ul>

### Related information:

Memory Kernel Services

Understanding Virtual Memory Manager Interfaces

## vm\_makep Kernel Service

### Purpose

Makes a page in client storage.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
int vm_makep ( vmid, pno)
vmid_t vmid;
int pno;
```

### Parameters

Item	Description
<i>vmid</i>	Specifies the ID of the virtual memory object.
<i>pno</i>	Specifies the page number in the virtual memory object.

### Description

The **vm\_makep** kernel service makes the page specified by the *pno* parameter addressable in the virtual memory object without requiring a page-in operation. The **vm\_makep** kernel service is restricted to client storage.

The page is not initialized to any particular value. It is assumed that the page is completely overwritten. If the page is already in memory, a value of 0, indicating a successful operation, is returned.

### Execution Environment

The **vm\_makep** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
EINVAL	Indicates a virtual memory object type or page number that is not valid.
EFBIG	Indicates that the page number exceeds the file-size limit.

### Related information:

Memory Kernel Services

Understanding Virtual Memory Manager Interfaces

## vm\_mem\_policy System Call Purpose

Allows callers to get or set their applications' default memory placement policies.

### Library

Standard C Library (**libc.a**)

### Syntax

```
#include <sys/rset.h>
```

```
#include <sys/vminfo.h>
```

```
int vm_mem_policy(int cmd, int *early_lru, int *policies, int num_policies)
```

### Description

The **vm\_mem\_policy** system call allows callers to get or set their applications' default memory placement policies for different types of memory.

Following are the different types of placement policies:

Item	Description
P_FIRST_TOUCH	Places the memory at the <b>MCM</b> where the application first referenced it. This is also achieved by setting the <b>MEMORY_AFFINITY</b> environment variable to <b>MCM</b> and benefit the applications with an identified home <b>MCM</b> to run on.
P_BALANCED	Uses the stripe memory in the application across all the system's <b>MCMs</b> . This benefits applications that do not identify a home <b>MCM</b> to run on, or on global memory objects that is accessed by many applications.
P_DEFAULT	Accepts the system's default policy for memory placement, which can be either the first touch or balanced policy, depending on the circumstances and the type of memory.

The **vm\_mem\_policy** system call allows the caller to get or set the **early\_lru** flag, which triggers the system to look for stealable pages immediately after a **P\_FIRST\_TOUCH** driven scan for local memory (the memory on the same **MCM** the application is running on) does not find any available pages.

The parameters *policies*, and *num\_policies* allow a caller to fine control over the default memory placement policies of different types of memory. The policy settings take effect on any new memory page the application creates after having called this function. The existing memory pages of the application retains their existing memory placement.

### Parameters

Item	Description
<i>cmd</i>	A command that is either <b>VM_SET_POLICY</b> or <b>VM_GET_POLICY</b> . The <b>VM_GET_POLICY</b> command copies the current policy setting into the buffers supplied by the caller, and does not change any of the process policies. The <b>VM_SET_POLICY</b> command reads input from the supplied buffers and changes the process policies accordingly.
<i>early_lru</i>	A pointer to an integer that indicates the state of the <i>early_lru</i> setting for first touch policy. Enabling <i>early_lru</i> causes memory to be paged out in order to fulfill a first-touch request for memory placement.  The possible values for <i>early_lru</i> are:  0        turn off <i>early_lru</i> . 1        turn on <i>early_lru</i> .  -1       do not modify <i>early_lru</i> setting for <b>VM_SET_POLICY</b> .
<i>policies</i>	A pointer to an array of policies for distinct types of memory. Each array element contains one of the policy types. The array element contains -1 to leave the policy unchanged for the corresponding memory type. The array must be declared with a length of <b>VM_NUM_POLICIES</b> . The list that follows enumerates the memory types whose policies can be changed in the form of constants. Enter the constant that is an array index into the policies array for the corresponding memory type.  <b>VM_POLICY_TEXT</b> policy for executable program text  <b>VM_POLICY_STACK</b> policy for program stack  <b>VM_POLICY_DATA</b> policy for program heap and private <b>mmap</b> data  <b>VM_POLICY_SHM_NAMED</b> policy for shared memory obtained via <b>shm_open()</b> or <b>shmget()</b> with a key  <b>VM_POLICY_SHM_ANON</b> policy for anonymous <b>mmap</b> memory, or shared memory obtained via <b>shmget()</b> with <b>IPC_PRIVATE</b> key  <b>VM_POLICY_MAPPED_FILE</b> policy for files mapped into the address space via <b>shmat()</b> or <b>mmap()</b>  <b>VM_POLICY_UNMAPPED_FILE</b> policy for open files that are not mapped
<i>num_policies</i>	Number of elements in the policies array. This value must be set to <b>VM_NUM_POLICIES</b> .

## vm\_mount Kernel Service

### Purpose

Adds a file system to the paging device table.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
int vm_mount ( type, ptr, nbufstr)
int type;
int (*ptr)();
int nbufstr;
```

### Parameters

Item	Description
<i>type</i>	Specifies the type of device. The <i>type</i> parameter must have a value of <b>D_REMOTE</b> .
<i>ptr</i>	Points to the file system's strategy routine.
<i>nbufstr</i>	Specifies the number of <b>buf</b> structures to use.

## Description

The **vm\_mount** kernel service allocates an entry in the paging device table for the file system. This service also allocates the number of **buf** structures specified by the *nbufstr* parameter for the calls to the strategy routine.

## Execution Environment

The **vm\_mount** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
ENOMEM	Indicates that there is no memory for the <b>buf</b> structures.
EINVAL	Indicates that the file system strategy pointer is already in the paging device table.

### Related reference:

“vm\_umount Kernel Service” on page 577

### Related information:

Memory Kernel Services

Understanding Virtual Memory Manager Interfaces

## vm\_moun~~te~~ Kernel Service

### Purpose

Adds a file system with a thread-level strategy routine to the paging device table.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
int vm_mounte ( in_dtype, in_devid, in_thrinfo)
int in_dtype;
dev_t in_devid;
struct thrpinfo * in_thrinfo;
```

### Parameters

Item	Description
<i>in_dtype</i>	Specifies the type of device. Supported device types are <b>D_REMOTE</b> , <b>D_LOGDEV</b> , <b>D_SERVER</b> , <b>D_LOCALCLIENT</b> . Other optional flags are detailed below.
<i>in_dev</i>	If the type is <b>D_LOGDEV</b> , specifies a <code>dev_t</code> object of the block device. If the type is <b>D_REMOTE</b> or <b>D_SERVER</b> , specifies a pointer to a strategy routine.
<i>in_thrinfop</i>	Pointer to a <b>thrpinfo</b> structure.

## Description

The **vm\_moun**te kernel service allocates an entry in the paging device table for the device specified. The **vm\_moun**te kernel service can also mount a client file system with a thread-level strategy routine. This is done by passing in the **D\_THRPGIO** and the **D\_ENHANCEDIO** flags.

## Flags

Name	Purpose
<b>D_ENHANCEDIO</b>	Indicates an enhanced I/O-aware file system.
<b>D_PREXLATE</b>	Enables pre-translation as the default for all but remote file systems.
<b>D_THRPGIO</b>	Indicates a thread-level strategy routine.

## Execution Environment

The **vm\_moun**te kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
ENOMEM	Indicates that there is no memory for the <b>buf</b> or the <b>thrpinfo</b> structure.
EINVAL	Indicates one of the following errors: <ul style="list-style-type: none"> <li>The file system strategy pointer is already in the paging device table, or in case of <b>D_SERVER</b>, a server is already defined.</li> <li>The <i>in_dtype</i> parameter is set to the <b>D_PAGING</b> or the <b>D_FILESYSTEM</b> value.</li> <li>The <b>thrpinfo</b> structure has not been initialized correctly.</li> <li>The <b>D_THRPGIO</b> flag has been set without the <b>D_ENHANCEDIO</b> flag.</li> </ul>

### Related reference:

“**vm\_umount** Kernel Service” on page 577

### Related information:

Memory Kernel Services

Understanding Virtual Memory Manager Interfaces

## vm\_move Kernel Service

### Purpose

Moves data between a virtual memory object and a buffer specified in the **uio** structure.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
#include <sys/uio.h>
```

```
int vm_move (vmid, offset, limit, rw, uio)
vmid_t vmid;
```

```

caddr_t  offset;
int     limit;
enum uio_rw  rw;
struct uio * uio;

```

## Parameters

Item	Description
<i>vmid</i>	Specifies the virtual memory object ID.
<i>offset</i>	Specifies the offset in the virtual memory object.
<i>limit</i>	Indicates the limit on the transfer length. If this parameter is negative or 0, no bytes are transferred.
<i>rw</i>	Specifies a read/write flag that gives the direction of the move. The possible values for this parameter ( <b>UIO_READ</b> , <b>UIO_WRITE</b> ) are defined in the <code>/usr/include/sys/uio.h</code> file.
<i>uio</i>	Points to the <b>uio</b> structure.

## Description

The **vm\_move** kernel service moves data between a virtual memory object and the buffer specified in a **uio** structure.

This service determines the virtual addressing required for the data movement according to the offset in the object.

The **vm\_move** kernel service is similar to the **uiomove** kernel service, but the address for the trusted buffer is specified by the *vmid* and *offset* parameters instead of as a **caddr\_t** address. The offset size is also limited to the size of a **caddr\_t** address since virtual memory objects must be smaller than this size.

**Note:** The **vm\_move** kernel service does not support use of cross-memory descriptors.

I/O errors for paging space and a lack of paging space are reported as signals.

## Execution Environment

The **vm\_move** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
EFAULT	Indicates a bad address.
ENOMEM	Indicates insufficient memory.
ENOSPC	Indicates insufficient disk space.
EIO	Indicates an I/O error.

Other file system-specific **errno** global variables are returned by the virtual file system involved in the move function.

### Related reference:

“uiomove Kernel Service” on page 517

### Related information:

Memory Kernel Services

Understanding Virtual Memory Manager Interfaces



## vm\_mvc Kernel Service

### Purpose

Reads or writes partial pages of files.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
int vm_mvc ( in_sid, in_pno, in_pgoffs, in_count, in_cmd, in_xmemdp, in_ptr)
vmid_t in_sid;
vpn_t in_pno;
int in_pgoffs;
int in_count;
int in_cmd;
struct xmem * in_xmemdp;
void * in_ptr;
```

### Parameters

Item	Description
<i>in_sid</i>	The primary memory object, m1.
<i>in_pno</i>	The m1 pno object. If it is a read operation, this parameter refers to the source. If not, it refers to a target.
<i>in_pgoffs</i>	The byte offset in the pno object.
<i>in_count</i>	The number of bytes to zero or copy in memory.
<i>in_cmd</i>	The reason for the function call. The possible values could be Zero, Zero(protect), read, or write.
<i>in_xmemdp</i>	The xmem descriptor for the second memory object, m2.
<i>in_ptr</i>	The byte offset in the xmem object.

### Description

The **vm\_mvc** kernel service is meant to be used by client file systems doing read or write operations to partial pages of files, where the file is denoted by the m1 object and the read or write buffer by the m2 object. Such cases arise on EOF handling, fragments, compression, and holes among other situations.

Given two memory object, m1 and m2, the **vm\_mvc** kernel service allows you to do one of the following operations:

- Zero out bytes on the m1 object (**VM\_MVC\_ZERO**).
- Zero out and protect the m1 object (**VM\_MVC\_PROTZERO**).
- Copy bytes from the m1 object to the m2 object (**VM\_MVC\_READ**).
- Copy bytes from the m2 object to the m1 object (**VM\_MVC\_WRITE**).

The first memory object, m1, is characterized by a *sid* parameter and a *pno* parameter. The second memory object, m2, is characterized by an xmem descriptor and a pointer for an offset. The second memory object is a user or kernel buffer.

**Note:** The second memory object must be pinned.

### Flags

<i>in_cmd</i>	<b>Purpose</b>
VM_MVC_ZERO	Zeros out the bytes on the m1 object.
VM_MVC_READ	Copies bytes from the m1 object to the m2 object.
VM_MVC_WRITE	Copies bytes from the m2 object to the m1 object.
VM_MVC_PROTZERO	Zeros out and protects the m1 object.

## Execution Environment

The `vm_mvc` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates that the I/O access operations completed successfully.
ENOENT	Indicates that the (sid, pno) set was not mapped to a real frame.
EINVAL	Indicates one of the following errors: <ul style="list-style-type: none"> <li>The m1 object crosses page boundary.</li> <li>The <i>in_cmd</i> parameter does not contain a valid command.</li> </ul>

### Related information:

Memory Kernel Services

Understanding Virtual Memory Manager Interfaces

## vm\_pattr System Call and kvm\_pattr Kernel Service Purpose

Queries or modifies virtual memory attributes.

## Library

Standard C Library (`libc.a`)

## Syntax

```
#include <sys/vmpattr.h>
```

```
int vm_pattr (
long cmd,
pid_t pid,
void * attr,
size_t attr_size );
```

```
int kvm_pattr (
long cmd,
pid_t pid,
void * attr,
size_t attr_size );
```

## Description

The `vm_pattr` system call queries or modifies memory attributes of the calling process's address space or that of another user process.

The `kvm_pattr` kernel service provides the same function to kernel subsystems (kernel extensions, kernel processes and so on) except that it cannot modify another kernel process' memory attributes.

## Parameters

**Item**  
*cmd*

### Description

The following commands can be passed in:

#### **VM\_PA\_SET\_PSIZE or VM\_PA\_GET\_PSIZE**

These commands set or retrieve the page size used for a specified memory range.

#### **VM\_PA\_GET\_RMUSAGE**

This command retrieves the amount of the real memory in bytes being used for a specified memory range.

#### **VM\_PA\_SET\_PSPA or VM\_PA\_GET\_PSPA**

These commands set or retrieve the page size promotion aggressiveness factor for a specified memory range.

#### **VM\_PA\_GET\_PSPA\_ALIGN**

This command retrieves the minimum memory alignment necessary for memory ranges specified to the **vm\_patr** kernel service when using the **VM\_PA\_SET\_PSPA** command.

#### **VM\_PA\_CHECK\_PSIZE**

This command reports if a specified page size can be used for a memory range.

#### **VM\_PA\_SET\_LSA\_POLICY**

This command allows the shared memory address space allocator to be tuned according to a process requirements. This command should be run before any shared memory regions are created.

#### **VM\_PA\_SET\_PSIZE\_EXTENDED or VM\_PA\_GET\_PSIZE\_EXTENDED**

These commands provide variable large page size segment support.

*pid*

Specifies the ID of the process whose memory attributes are to be queried or modified. A value of -1 specifies the calling process. The root user can specify any process ID, but other users can only specify processes they own (that is, the target process's user ID must match the calling process's user ID).

The **vm\_patr** system call is only supported on user processes while the **kvm\_patr** kernel service can target user processes or its own kernel process (for example, *pid* = -1).

*attr*

A pointer to a structure describing the effective address range for the memory being queried or modified and additional data depending on the command.

The range is specified through the following **vm\_pa\_range** structure:

```
struct vm_pa_range
{
    ptr64_t  rng_start;
    size64_t rng_size;
};
```

The range specified must be in the target process's address space and must correspond to one of these process areas:

- Main program data (initialized, bss, or heap).
- Shared library data or private module load area data.
- Privately loaded text.
- Initial thread stack area.
- Anonymous shared memory (System V shared memory, extended System V shared through EXTSHM, and POSIX real-time shared memory). The target process must have write access to the memory in order to change the attributes of the shared memory range.
- Anonymous mmap memory.

If the memory range specified includes shared memory or mmap memory, the calling process must have write access to the memory according to the shared memory descriptor or mapping attributes in order to change the attributes of the range. The range can have additional restrictions based on the following commands.

**Item***attr* (continued)**Description**

The structure specified through the *attr* parameter must be a pointer to one of the following structures:

**VM\_PA\_SET\_PSIZE or VM\_PA\_GET\_PSIZE**

These commands take a pointer to the following structure:

```
struct vm_pa_psize
{
    struct vm_pa_range pa_range;
    psize_t           pa_psize;
};
```

For the **VM\_PA\_SET\_PSIZE** command, the *pa\_psize* parameter is the page size (in bytes) to use for the given range. This is an advisory setting that might or might not be used at the operating system's discretion. This must be a valid page size between the minimum and maximum page sizes of all segments in the range. Additionally, the range must start and end on a multiple of the specified page size. If an error occurs during the processing of this command, any successfully altered page size settings can remain set.

For the **VM\_PA\_GET\_PSIZE** command, the page size (in bytes) backing the specified memory range is returned in the *pa\_psize* parameter. The range must start and end on a multiple of the smallest page size supported as reported by the `sysconf(SC_PAGE_SIZE)` subroutine. If the range is using multiple page sizes, the smallest page size in the range is reported. Unlike the **VM\_PAGE\_INFO** command of the `vmgetinfo` subroutine that reports the segment's base page size, the page size reported by the **VM\_PA\_GET\_PSIZE** command is the actual page size being used at the time the `vm_patr` system call was called. The page size reported is transient because the operating system can change the backing page size at any time. Therefore, the page size reported must be for informational purposes only.

**VM\_PA\_SET\_PSIZE\_EXTENDED**

This command takes a pointer to the following structure:

```
struct vm_pa_psize_extended
{
    struct vm_pa_range pa_range;
    psize_t           pa_psize;
    size_t            pa_info_size;
    uint64_t          *pa_info;
};
```

This command is essentially the same as **VM\_PA\_SET\_PSIZE** except that *pa\_psize* must be 16 MB and, if not NULL, *pa\_info* can be used to pass additional information specifying one or more affinity domains.

The info passed by the parameter is advisory request, and the system might choose to ignore it.

The *pa\_range* is scanned for subregions that begin and end on a 16 MB boundary, are fully backed with 4 KB or 64 KB pages, and have uniform page attributes. The page attributes include read or write page protection, storage key protection, and no-execute protection.

The data in qualifying 16 MB subregions is collocated to a 16 MB contiguous block of physical memory, and it uses 16 MB hardware translations.

If the *pa\_info* pointer is NULL, the memory for collocation is allocated from any memory SRAD, affinity domain chosen by the operating system.

If parameter value is not NULL, *pa\_info* must point to an *rsethandle\_t* that describes a set of affinity domains from which the physical memory for the collocation must be allocated. The object should be allocated by a call to `rs_alloc(RS_EMPTY)`. It must then be initialized with one `rs_op(RS_ADDRESOURCE, ..., R_MEMPS, srads#)` call per affinity domain being requested.

This command can potentially affect system performance and is not generally recommended; therefore, this command requires you to have either the `CAP_BYPASS_RAC_VMM` and `CAP_PROPAGATE` capabilities or root authority.

**VM\_PA\_GET\_PSIZE\_EXTENDED**

This command is essentially the same as the **VM\_PA\_GET\_PSIZE** command except that it can also return the 64 KB and 16 MB subregions that are using an hardware translation page size different from the underlying segments default page size.

**Item**  
*attr* (continued)

**Description**

If the *pa\_info* field is NULL, this command is identical to the **VM\_PA\_GET\_PSIZE** command.

The *pa\_info* field should point to an array containing two 64 bit integers. The *pa\_info\_size* field should be set to the size of the array.

In the first 64-bit integer, this command reports the number of 64-KB sized and aligned subregions in the specified *pa\_range* range that consist of 16 contiguous 4-KB pages that are promoted to using a 64-KB page size hardware translation. In the second 64-bit integer, this command reports the number of 16 MB sized and aligned subregions in the specified range that consist of either 4096 4 KB or 256 64-KB contiguous pages that are promoted to using a 16-MB page size hardware translation.

The *pa\_psize* field reports the smallest page size found for the specified range.

The information reported is transient because the operating system can change the backing page size at any time. Therefore, the page size reported must be for informational purposes only.

*attr* (continued)

**VM\_PA\_GET\_RMUSAGE**

This command takes a pointer to the following structure:

```
struct vm_pa_rmusage
{
    struct vm_pa_range pa_range;
    size64_t          pa_rbytes;
};
```

This command reports the amount of real memory (in bytes) used for the given range in the *pa\_rbytes* field. This can help an application decide whether it needs to use a large page size for a specific range based on how much real memory the range is using. For example, if a 64KB range is only using 4KB of real memory, then it does not make sense to try to use a 64KB page size for that range. But if it is using all 64KB or some large percentage of it, then the application might decide to use a 64KB page size. The range specified for this command has no alignment requirements for this command, and the command includes only those bytes in the range that are using real memory.

**VM\_PA\_SET\_PSPA or VM\_PA\_GET\_PSPA**

These commands take a pointer to the following structure:

```
struct vm_pa_pspa
{
    struct vm_pa_range pa_range;
    int pa_pspa;
};
```

The **VM\_PA\_SET\_PSPA** command can set the page size promotion aggressiveness for the specified range. The *pa\_pspa* setting is in the same units as the `vmm_default_pspa` vmo tunable. This setting is the inverse of the real memory occupancy threshold needed to promote to a large page size and ranges from -1 to 100. The value of -1 indicates that no page promotion can occur regardless of the occupancy of the memory range. A value of 0 indicates a page size promotion can only be done when the memory range is fully occupied. A value of 100 indicates a page promotion must be done at the first reference to the memory range.

This setting is only supported at a segment granularity, so the range must start and end on a segment boundary. The alignment requirement for the range can be found using the **VM\_PA\_GET\_PSPA\_ALIGN** command with the `vm_pattr` system call.

If an error occurs during the processing of the **VM\_PA\_SET\_PSPA** command, the `vm_pattr` system call can return after altering the page size promotion thresholds for part of the specified range.

The **VM\_PA\_GET\_PSPA** command retrieves the page size promotion aggressiveness factor for the specified range. If the range spans multiple segments consisting of different page promotion thresholds, the *pa\_pspa* field is updated with the least aggressive PSPA setting (the smallest PSPA setting across all of the segments).

The PSPA commands are not supported on `mmap` or `EXTSHM` memory ranges.

**Item****Description***attr* (continued)**VM\_PA\_GET\_PSPA\_ALIGN**

This command takes a pointer to the following structure:

```
struct vm_pa_pspa_align
{
    struct vm_pa_range pa_range;
    size64_t pa_pspa_align;
};
```

The **VM\_PA\_GET\_PSPA\_ALIGN** command returns the minimum memory alignment requirements of a memory range for the **VM\_PA\_SET\_PSPA** command in the *pa\_pspa\_align* field based on what segments are contained in the specified memory range. If a memory range spans segments with different alignment requirements, this command returns the largest of the alignment requirements.

The alignment requirements for the **VM\_PA\_SET\_PSPA** command are as follows:

*attr* (continued)**VM\_PA\_SET\_LSA\_POLICY**

This command takes a pointer to the following structure:

```
struct vm_pa_lsa_options
{
    u_int64_t setting;
    size64_t value;
};
```

The following settings are allowed:

**VM\_PA\_SHM\_1TB\_SHARED**

This setting controls the threshold of the number of 256 MB segments required before a SHM object is considered big enough to be placed in its own 1 TB region to be promoted to the large alias segments. Values can range from 0 to 4 KB.

**VM\_PA\_SHM\_1TB\_UNSHARED**

This setting controls the threshold of the number of 256 MB segments required before a group of SHM object packed in a 1 TB aligned group is promoted to the large alias segments. Values can range from 0 to 4 KB.

*attr* (continued)**Process's Memory Area Minimum Alignment**

Main process data 256 MB

Process stack 256 MB

Shared Library data 256 MB

Privately loaded module data 256 MB

Privately loaded module text 256 MB

POSIX Real-Time Shared Memory 256 MB

Anonymous MMAP 256 MB

Anonymous Extended System V Shared memory 256 MB

Anonymous System V Shared memory with page sizes less than or equal to 256 MB 256 MB

Anonymous System V shared memory backed with 16 GB page size 1 TB

**Item**  
*attr* (continued)

**Description**

**VM\_PA\_CHECK\_PSIZE**

This command takes a pointer to the following structure:

```
struct vm_pa_psize_check
{
    struct vm_pa_range pa_range;
    psize_t           pa_psize;
    int               pa_reason;
};
```

The **VM\_PA\_CHECK\_PSIZE** command determines if a specific page size is allowed by the **VM\_PA\_SET\_PSIZE** command for a specified memory range. The **VM\_PA\_CHECK\_PSIZE** command can be used when the application wants more detailed information about why a **VM\_PA\_SET\_PSIZE** operation fails, or to check if a **VM\_PA\_SET\_PSIZE** operation will successfully modify the page size for the range specified.

This command must be used on a memory range that spans a single page and is aligned to the page size specified by the *pa\_psize* parameter. If the page size can be used for that range, the *pa\_reason* parameter is set to 0. Otherwise, it is set to a reason code defined in the **vmpattr.h** header file.

**VMPATTR\_SET\_PSIZE\_VALID** The specified page size can be used for the specified range.

**VMPATTR\_INVALID\_MPSS\_PSIZE** The specified page size is not supported in mixed page size segments.

**VMPATTR\_NON\_MPSS\_SEGMENT** The address range specified is from a segment that does not support mixed page sizes.

**VMPATTR\_NON\_MPSS\_PAGE** The size of the target page cannot be modified. For example, this reason code can be returned when trying to set an address range to a 64 KB page size if a portion of the range has page protection settings that do not match the rest of the range.

**VMPATTR\_RDONLY\_MEM** The target range cannot be modified because the caller does not have write access to the memory specified.

**VMPATTR\_PAGE\_ATTRIBUTES** The address range specified does not have uniform page attributes.

**VMPATTR\_NOT\_FULLY\_POPULATED** The address range specified does not fully reside in memory.

**VMPATTR\_PHYSICAL\_ATTACHMENTS** The address range specified has memory affinity attachments that specify more than one affinity domain.

**VMPATTR\_MEMORY\_TYPE\_UNSUPPORTED** The address range contains a memory object that does not support the requested page size in a mixed page size segment.

*attr\_size*

The *attr\_size* parameter must be the size of the structure needed, or greater for the specified command.

## Return Values

When successful, these commands return 0. Otherwise, they return -1 and set the `errno` global variable to indicate the error.

## Error Codes

Item	Description
EPERM	The calling process does not have the appropriate privilege to perform the requested operation.
ESRCH	The target process does not exist or is not in a valid state.
ENOMEM	The range specified contains a hole. A hole is any part of the target process's address space that is not backed by a virtual memory segment or is outside of the valid range of the virtual memory segment specified.
ENOTSUP	Any of the following situations can cause the ENOTSUP error: <ul style="list-style-type: none"> <li>• The target process is a kernel process other than the calling process.</li> <li>• The command specified was the <code>VM_PA_SET_PSIZE</code> command and the page size specified is not supported for multiple page size segments.</li> <li>• The command specified was either the <code>VM_PA_GET_PSPA</code> or the <code>VM_PA_SET_PSPA</code> command and the specified memory range includes mmap or EXTSHM segment(s).</li> </ul>
EINVAL	Any of the following situations can cause the EINVAL error: <ul style="list-style-type: none"> <li>• The <code>attr_size</code> parameter specified is less than the size of the structure needed for this command.</li> <li>• The range specified is outside the process's address space (for example, global kernel memory).</li> <li>• The command specified was the <code>VM_PA_SET_PSIZE</code> command and the page size specified was not a valid page size supported by the system.</li> <li>• The command specified was the <code>VM_PA_SET_PSPA</code> command and the address range specified was not aligned to the segment size backing the range.</li> <li>• The command specified was the <code>VM_PA_SET_PSPA</code> command and the page promotion aggressiveness factor specified was not valid.</li> <li>• The command specified was the <code>VM_PA_CHECK_PSIZE</code> command and the address range specified was not aligned to the page size specified.</li> </ul>
ENOMEM	The command specified was <code>VM_PA_SET_PSIZE_EXTENDED</code> , and the system was unable to allocate memory from the set of affinity domains specified by the <code>pa_info</code> object or the entire set of system affinity domains without potentially causing a performance degradation.
EFAULT	The command specified was either <code>VM_PA_SET_PSIZE_EXTENDED</code> , or <code>VM_PA_GET_PSIZE_EXTENDED</code> and the <code>pa_info</code> address is not valid and is not NULL.
EINVAL	The command specified was either <code>VM_PA_SET_PSIZE_EXTENDED</code> , or <code>VM_PA_GET_PSIZE_EXTENDED</code> and the <code>pa_info</code> field is not-NULL, but the <code>pa_info_size</code> field is 0.
ENODEV	The command specified was <code>VM_PA_SET_PSIZE_EXTENDED</code> , and an invalid <code>sraddr</code> was specified in <code>pa_info</code> .

#### Related information:

Dynamic variable page size support

### vm\_protect\_kkey Kernel Service Purpose

Sets kernel-key on a kernel address range.

#### Syntax

```
#include <sys/types.h>
#include <sys/skeys.h>
#include <sys/vmuser.h>
```

```
kernno_t vm_protect_kkey (eaddr, nbytes, kkey, flags)
void * eaddr;
size_t nbytes;
kkey_t kkey;
unsigned long flags;
```



## Parameters

Item	Description
<i>eaddr</i>	Starting address to protect.
<i>nbytes</i>	Number of bytes to protect.
<i>kkey</i>	Kernel-key value to set on memory.
<i>flags</i>	Defined flag value is: <ul style="list-style-type: none"><li>• <b>VMPK_NO_CHECK_AUTHORITY</b> – This flag indicates that extended authority checking will not be performed.</li></ul>

## Description

The `vm_protect_kkey()` kernel service is used to alter the kernel-key associated with a virtual memory range. If set, any code that references the memory needs to include the kernel-key in their active keyset. The kernel-key is set for all pages in the effective address range specified by *eaddr* to *eaddr + nbytes - 1*. If the address range does not specify a page-aligned area consisting of an integral number of full pages, an error will be returned.

By default, an authority check is performed when altering storage-keys. This check requires that the `vm_protect_kkey()` caller has write access to the pages' current kernel-key(s). This authority checking can be overridden by setting the **VMPK\_NO\_CHECK\_AUTHORITY** value, but this is not recommended since the check can protect against some programming errors.

## Execution Environment

The `vm_protect_kkey` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Successful.
<b>EINVAL_VM_PROTECT_KKEY</b>	Invalid parameter or execution environment.
<b>EINVAL_VM_PROTECT_KKEY_PPAGE</b>	Request includes a partial page.
<b>EFAULT_VM_PROTECT_KKEY</b>	Invalid address range.
<b>EPERM_VM_PROTECT_KKEY</b>	Insufficient authority to perform the operation.

If the `vm_protect_kkey()` kernel service is unsuccessful because of a condition other than that specified by the **EINVAL\_VM\_PROTECT\_KKEY** error code, the kernel-key for some pages in the (*eaddr*, *eaddr + nbytes - 1*) range might have been changed.

### Related reference:

“`vm_setseg_kkey` Kernel Service” on page 570

## vm\_protectp Kernel Service

### Purpose

Sets the page protection key for a page range.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
int vm_protectp ( vmid, pfirst, npages, key)
vmid_t vmid;
```

```
int pfirst;
int npages;
int key;
```

## Description

The **vm\_protectp** kernel service is called to set the storage protect key for a given page range. The *key* parameter specifies the value to which the page protection key is set. The protection key is set for all pages touched by the specified page range that are resident in memory. The **vm\_protectp** kernel service applies only to client storage.

If a page is not in memory, no state information is saved from a particular call to the **vm\_protectp** service. If the page is later paged-in, it receives the default page protection key.

**Note:** The **vm\_protectp** subroutine is not supported for use on large pages.

## Parameters

Item	Description
<i>vmid</i>	Specifies the identifier for the virtual memory object for which the page protection key is to be set.
<i>pfirst</i>	Specifies the first page number in the designated page range.
<i>npages</i>	Specifies the number of pages in the designated page range.
<i>key</i>	Specifies the value to be used in setting the page protection key for the designated page range.

## Execution Environment

The **vm\_protectp** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
EINVAL	Indicates one of the following errors: <ul style="list-style-type: none"><li>• Invalid virtual memory object ID.</li><li>• The starting page in the designated page range is negative.</li><li>• The number of pages in the page range is negative.</li><li>• The designated page range exceeds the size of virtual memory object.</li><li>• The target page range does not exist.</li><li>• One or more large pages lie in the target page range.</li></ul>

### Related information:

Memory Kernel Services

Understanding Virtual Memory Manager Interfaces

## vm\_qmodify Kernel Service

### Purpose

Determines whether a mapped file has been changed.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
int vm_qmodify ( vmid)
vmid_t vmid;
```

## Parameter

Item	Description
<i>vmid</i>	Specifies the ID of the virtual memory object to check.

## Description

The **vm\_qmodify** kernel service performs two tests to determine if a mapped file has been changed:

- The **vm\_qmodify** kernel service first checks the virtual memory object modified bit, which is set whenever a page is written out.
- If the modified bit is 0, the list of page frames holding pages for this virtual memory object are examined to see if any page frame has been modified.

If both tests are false, the **vm\_qmodify** kernel service returns a value of False. Otherwise, this service returns a value of True.

If the virtual memory object modified bit was set, it is reset to 0. The page frame modified bits are not changed.

## Execution Environment

The **vm\_qmodify** kernel service can be called from the process environment only.

## Return Values

Item	Description
FALSE	Indicates that the virtual memory object has not been modified.
TRUE	Indicates that the virtual memory object has been modified.

## Related information:

Memory Kernel Services

Understanding Virtual Memory Manager Interfaces

## vm\_qppages Kernel Service

### Purpose

Returns the number of in-memory page frames associated with the virtual memory object.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
vpn_t vm_qppages ( sid)
vmid_t sid;
```

### Parameters

Item	Description
<i>sid</i>	Identifies the base segment.

## Description

The **vm\_qpages** kernel service routine returns the number of page frames associated with the virtual memory object with the *sid* parameter specified.

## Execution Environment

The **vm\_qpages** kernel service can be called from the process environment only.

This function can be run for persistent, client, and working storage segments.

## Return Values

Item	Description
<b>npages</b>	The number of page frames.
<b>-1</b>	Indicates an invalid <i>sid</i> parameter.

## Related information:

Memory Kernel Services

Understanding Virtual Memory Manager Interfaces

## vm\_readp Kernel Service

### Purpose

Initiates asynchronous page-in for the range of pages specified.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
int vm_readp ( sid, pfirst, npages, flags)
vmid_t sid;
vpn_t pfirst;
vpn_t npages;
int flags;
```

### Parameters

Item	Description
<i>sid</i>	Identifies the base segment.
<i>pfirst</i>	The first page number within the range.
<i>npages</i>	The number of pages to read starting from the <i>pfirst</i> value. All pages must be in the same segment, unless the <b>V_READMAKE</b> flag is used.
<i>flags</i>	Flags used by the function.

### Description

- | The **vm\_readp** kernel service routine starts the process of paging within the range of specified pages. This call is strictly advisory and might return without performing any operations.
- | The *flags* parameter is optional and accepts the following values:

- | **V\_IOWAIT**  
|       Instructs the **vm\_readp** kernel service to wait for any page I/O requests to complete, within the  
|       range of specified pages, before initiating the read operation.
- | **V\_READMAKE**  
|       Instructs the **vm\_readp** kernel service to create the segments within the range of the **vm\_readp**  
|       operation.

## Execution Environment

The **vm\_readp** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates that the I/O access operations completed successfully.
EINVAL	Indicates one of the following errors: <ul style="list-style-type: none"> <li>• <i>ppfirst</i> = 0 and <i>npages</i> = 0.</li> <li>• <i>ppfirst</i> &lt; 0.</li> <li>• <i>npages</i> &lt; 0.</li> <li>• Page interval &gt; Maximum file size.</li> <li>• The <i>sid</i> parameter is not valid.</li> <li>• Not a file or persistent storage segment.</li> </ul>

### Related reference:

“vm\_writep Kernel Service” on page 579

### Related information:

Memory Kernel Services

Understanding Virtual Memory Manager Interfaces

## vm\_release Kernel Service

**Note:** The **vm\_release** subroutine is not supported for use on large pages.

### Purpose

Releases virtual memory resources for the specified address range.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
int vm_release ( vaddr, nbytes)
caddr_t vaddr;
int nbytes;
```

### Description

The **vm\_release** kernel service releases pages that intersect the specified address range from the *vaddr* parameter to the *vaddr* parameter plus the number of bytes specified by the *nbytes* parameter. The value in the *nbytes* parameter must be nonnegative and the caller must have write access to the pages specified by the address range.

Each page that intersects the byte range is logically reset to 0, and any page frame is discarded. A page frame in I/O state is marked for discard at I/O completion. That is, the page frame is placed on the free list when the I/O operation completes.

**Note:** All of the pages to be released must be in the same virtual memory object.

## Parameters

Item	Description
<i>vaddr</i>	Specifies the address of the first byte in the address range to be released.
<i>nbytes</i>	Specifies the number of bytes to be released.

## Execution Environment

The `vm_release` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates successful completion.
EACCES	Indicates that the caller does not have write access to the specified pages.
EINVAL	Indicates one of the following errors: <ul style="list-style-type: none"><li>• The specified region is not mapped.</li><li>• The specified region is an I/O region.</li><li>• The length specified in the <i>nbytes</i> parameter is negative.</li><li>• The specified address range crosses a virtual memory object boundary.</li><li>• One or more large pages lie in the target page range.</li></ul>

### Related reference:

“`vm_releasep` Kernel Service”

### Related information:

Memory Kernel Services

Understanding Virtual Memory Manager Interfaces

## `vm_releasep` Kernel Service

### Purpose

Releases virtual memory resources for the specified page range.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
int vm_releasep ( vmid, pfirst, npages)
vmid_t vmid;
int pfirst;
int npages;
```

### Description

The `vm_releasep` kernel service releases pages for the specified page range in the virtual memory object. The values in the *pfirst* and *npages* parameters must be nonnegative.

Each page of the virtual memory object that intersects the page range (*pfirst*, *pfirst* + *npages* -1) is logically reset to 0, and any page frame is discarded. A page frame in the I/O state is marked for discard at I/O completion.

For working storage, paging-space disk blocks are freed and the storage-protect key is reset to the default value.

**Note:** All of the pages to be released must be in the same virtual memory object.

**Note:** The `vm_releasep` subroutine is not supported for use on large pages.

## Parameters

Item	Description
<i>vmid</i>	Specifies the virtual memory object identifier.
<i>pfirst</i>	Specifies the first page number in the specified page range.
<i>npages</i>	Specifies the number of pages in the specified page range.

## Execution Environment

The `vm_releasep` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
EINVAL	Indicates one of the following errors: <ul style="list-style-type: none"><li>• An invalid virtual memory object ID.</li><li>• The starting page is negative.</li><li>• Number of pages is negative.</li><li>• Page range crosses a virtual memory object boundary.</li><li>• One or more large pages lie in the target page range.</li></ul>

### Related reference:

“`vm_release` Kernel Service” on page 565

### Related information:

Memory Kernel Services

Understanding Virtual Memory Manager Interfaces

## `vm_segmap` Kernel Service

### Purpose

Creates the segments associated with a range of bytes in a file and attaches them to the kernel's address space.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
int vm_segmap ( basesid, pfirst, flags, basepp)
vmid_t basesid;
vpn_t pfirst;
uint flags;
caddr_t * basepp;
```

## Parameters

Item	Description
<i>basesid</i>	Identifies the base segment.
<i>pfirst</i>	The first page number within the range. The <i>pfirst</i> parameter is non-negative.
<i>flags</i>	Optional flags passed in by the user. .
<i>basepp</i>	The offset of the object to be attached.

## Description

The **vm\_segmap** kernel service routine creates segments associated with a range of bytes in a file. Afterwards, it uses the **vm\_att** kernel service to map the specified virtual memory object to a region in the virtual address space and returns the effective address of that object in the *basepp* parameter.

## Execution Environment

The **vm\_segmap** kernel service can be called from either the process or interrupt environment.

## Return Values

Item	Description
<b>caddr_t</b>	The effective address of the attached object.
<b>EINVAL</b>	Indicates one of the following errors: <ul style="list-style-type: none"><li>• <i>pfirst</i> &lt; 0.</li><li>• Invalid <i>sid</i> parameter.</li></ul>
<b>EFBIG</b>	Indicates the range of values is too large to create.

### Related information:

Memory Kernel Services

Understanding Virtual Memory Manager Interfaces

## vm\_setdevid Kernel Service

### Purpose

Modifies the paging device table entry for a virtual memory object.

### Syntax

```
#include <sys/types.h>
#include <sys/kernno.h>
#include <sys/vmuser.h>
```

```
kernno_t vm_setdevid ( vmid, type, ptr, flags)
vmid_t vmid;
int type;
int (*ptr)();
unsigned long flags;
```

### Parameters



Item	Description
<i>vmid</i>	Specifies the identifier for the virtual memory object for which the paging device table entry is to be set.
<i>type</i>	Specifies the type of device. The <i>type</i> parameter must have a value of D_REMOTE.
<i>ptr</i>	Points to the strategy routine of the file system.
<i>flags</i>	Reserved. You must set the <i>flags</i> parameter to zero.

## Description

The **vm\_setdevid** kernel service binds the paging device table entry associated with the file system strategy routine *ptr*, to the virtual memory object *vmid*. The paging device table entry must have already been mounted as type D\_REMOTE through a prior **vm\_mount** kernel service call.

After the file system has called the **vm\_setdevid** kernel service on a given virtual memory object, subsequent paging I/O will be performed to or from the newly specified paging device table. Any outstanding I/O's to the paging device table formerly associated with the virtual memory object, remain queued, and will complete asynchronously. After they complete, subsequent paging I/O to those file pages will be performed to or from the newly specified paging device table.

The paging device table entry currently associated with the *vmid* object, on input to this call, must be valid and of type D\_REMOTE. Any flags specified when the **vm\_mount** kernel service gets called must match exactly any flags specified when the **vm\_mount** kernel service gets called for the new paging device table entry.

## Execution Environment

The **vm\_setdevid** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
EINVAL_VM_SETDEVID1	Indicates that the <i>vmid</i> value is not a client segment, or the input type does not have the value of D_REMOTE.
ENODEV_VM_SETDEVID2	Indicates that a file system with the strategy routine designated by the <i>ptr</i> parameter is not in the paging device table.
EINVAL_VM_SETDEVID3	Indicates that the new paging device table entry is not D_REMOTE or is not valid.
EINVAL_VM_SETDEVID4	Indicates that the paging device table entry currently associated with the <i>vmid</i> object is not D_REMOTE or is not valid.
EINVAL_VM_SETDEVID5	Indicates that the <b>vm_mount</b> flags for the current and new paging device table entries differ.
EINVAL_VM_SETDEVID6	Indicates that this was called at interrupt level.
EINVAL_VM_SETDEVID7	Indicates that the input flags was nonzero.
EINVAL_VM_SETDEVID8	Indicates that the input <i>vmid</i> value is not valid.

## Related Information

The **vm\_mount** kernel service, **vm\_umount** kernel service.

Memory Kernel Services and Understanding Virtual Memory Manager Interfaces in *Kernel Extensions and Device Support Programming Concepts*.

### Related reference:

“vm\_mount Kernel Service” on page 549

“vm\_umount Kernel Service” on page 577

### Related information:

Understanding Virtual Memory Manager Interfaces

## vm\_setseg\_kkey Kernel Service

### Purpose

Sets the default kernel-key for a segment.

### Syntax

```
#include <sys/types.h>
#include <sys/kernno.h>
#include <sys/vmuser.h>
```

```
kernno_t vm_setseg_kkey (vmid, kkey)
vmid_t vmid;
kkey_t kkey;
```

### Parameters

Item	Description
<i>vmid</i>	Virtual memory object to act on.
<i>kkey</i>	New kernel key for the virtual memory object.

### Description

The **vm\_setseg\_kkey** kernel service alters the default kernel-key for newly allocated pages in a segment. The kernel-key values for any existing pages in the segment are left unchanged.

### Execution Environment

The **vm\_setseg\_kkey** kernel service can be called from the process environment only.

### Return Values

Item	Description
0	Successful.
EINVAL_VM_SETSEG_KKEY	Invalid parameter or execution environment.

### Related reference:

“vm\_protect\_kkey Kernel Service” on page 560

## vm\_thrpgio\_pop Kernel Service

### Purpose

Retrieves the latest context information.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
void vm_thrpgio_pop ( in_ctxp)
ut_pgio_context_t * in_ctxp;
```

### Parameters

Item	Description
<i>in_ctxp</i>	The context structure used by the function.

## Description

The **vm\_thrpgio\_pop** kernel service enables a client file system with a thread-level strategy routine to copy information from a context structure to the current thread. Afterwards, it makes the current thread point to the next context.

This service must be called if a client file system using a thread-level strategy routine has re-entered the Virtual Memory Manager and wishes to return to its strategy routine. This service restores the context that was saved using the **vm\_thrpgio\_push** kernel service.

## Execution Environment

The **vm\_thrpgio\_pop** kernel service can only be used by client file systems using a thread-level strategy routine.

## Return Values

The **vm\_thrpgio\_pop** kernel service has no return values.

### Related reference:

“vm\_thrpgio\_push Kernel Service”

### Related information:

Memory Kernel Services

Understanding Virtual Memory Manager Interfaces

## vm\_thrpgio\_push Kernel Service Purpose

Saves some context information of the current thread.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
void vm_thrpgio_push ( in_ctxp)
ut_pgio_context_t * in_ctxp;
```

## Parameters

Item	Description
<i>in_ctxp</i>	The context structure used by the function.

## Description

The **vm\_thrpgio\_push** kernel service enables a client file system with a thread-level strategy routine to save information about the current thread to a linked list. The linked list is a Last-In-First-Out (LIFO) (stack) data structure, and is pointed to by the thread.

This service must be called if a client file system using a thread-level strategy routine has had its strategy routine invoked and wishes to re-enter the Virtual Memory Manager. This could involve a page fault on one of its client segments, or the use of one of the Virtual Memory Manager (VMM) services that operates on client segments.

The `vm_thrpgio_pop` kernel service must be invoked when all such Virtual Memory Manager callbacks are complete.

## Execution Environment

The `vm_thrpgio_push` kernel service can only be used by client file systems using a thread-level strategy routine.

## Return Values

The `vm_thrpgio_push` kernel service has no return values.

### Related reference:

“`vm_thrpgio_pop` Kernel Service” on page 570

### Related information:

Memory Kernel Services

Understanding Virtual Memory Manager Interfaces

## vms\_create Kernel Service

### Purpose

Creates a virtual memory object of the specified type, size, and limits.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
int vms_create (vmid, type, devgno, size, uplim, downlim)
vmid_t * vmid;
int type;
dev_t devgno;
int size;
int uplim;
int downlim;
```

### Parameters

Item	Description
<i>vmid</i>	Points to the variable in which the virtual memory object identifier is to be stored.
<i>type</i>	Specifies the virtual memory object type and options as an OR of bits. The <i>type</i> parameter must have the value of <code>V_CLIENT</code> . The <code>V_INTRSEG</code> flag specifies if the process can be interrupted from a page wait on this object.
<i>devgno</i>	Specifies the address of the g-node for client storage. If the <i>type</i> parameter has the value of <code>V_CLIENT</code> , the third argument is a g-node <i>ptr</i> argument, otherwise, it is a <i>devgno</i> argument.
<i>size</i>	Specifies the current size of the file (in bytes). This can be any valid file size. If the <code>V_LARGE</code> is specified, it is interpreted as number of pages.
<i>uplim</i>	Ignored. The enforcement of file size limits is done by comparing with the <code>u_limit</code> value in the <code>u</code> block.
<i>downlim</i>	Ignored.

## Description

The `vms_create` kernel service creates a virtual memory object. The resulting virtual memory object identifier is passed back by reference in the `vmid` parameter.

The `size` parameter is used to determine the size in units of bytes of the virtual memory object to be created. This parameter sets an internal variable that determines the virtual memory range to be processed when the virtual memory object is deleted.

An entry for the file system is required in the paging device table when the `vms_create` kernel service is called.

## Execution Environment

The `vms_create` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
ENOMEM	Indicates that no space is available for the virtual memory object.
ENODEV	Indicates no entry for the file system in the paging device table.
EINVAL	Indicates incompatible or bad parameters.

### Related reference:

“`vms_delete` Kernel Service”

### Related information:

Memory Kernel Services

Understanding Virtual Memory Manager Interfaces

## `vms_delete` Kernel Service

### Purpose

Deletes a virtual memory object.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
int vms_delete ( vmid)
vmid_t vmid;
```

### Parameter

Item	Description
<i>vmid</i>	Specifies the ID of the virtual memory object to be deleted.

## Description

The **vms\_delete** kernel service deallocates the temporary resources held by the virtual memory object specified by the *vmid* parameter and then frees the control block. This delete operation can complete asynchronously, but the caller receives a synchronous return code indicating success or failure.

## Releasing Resources

The completion of the delete operation can be delayed if paging I/O is still occurring for pages attached to the object. All page frames not in the I/O state are released.

If there are page frames in the I/O state, they are marked for discard at I/O completion and the virtual memory object is placed in the iodelete state. When an I/O completion occurs for the last page attached to a virtual memory object in the iodelete state, the virtual memory object is placed on the free list.

## Execution Environment

The **vms\_delete** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
EINVAL	Indicates that the <i>vmid</i> parameter is not valid.

## Related reference:

“vms\_create Kernel Service” on page 572

## Related information:

Memory Kernel Services

Understanding Virtual Memory Manager Interfaces

## vms\_iowait, vms\_iowaitf Kernel Services Purpose

Waits for the completion of all page-out operations for pages in the virtual memory object.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
int vms_iowait ( vmid)
vmid_t vmid;
int vms_iowaitf ( vmid, flags)
vmid_t vmid;
int flags;
```

## Parameter

Item	Description
<i>vmid</i>	Identifies the virtual memory object for which to wait.
<i>flags</i>	Optional flags passed in by the user.

## Description

The **vms\_iowait** kernel service performs two tasks. First, it determines the I/O level at which all currently scheduled page-outs are complete for the virtual memory object specified by the *vmid* parameter. Then, the **vms\_iowait** service places the current process in a wait state until this I/O level has been reached.

The I/O level value is a count of page-out operations kept for each virtual memory object.

The I/O level accounts for out-of-order processing by not incrementing the I/O level for new page-out requests until all previous requests are complete. Because of this, processes waiting on different I/O levels can be awakened after a single page-out operation completes.

If the caller holds the kernel lock, the **vms\_iowait** service releases the kernel lock before waiting and reacquires it afterwards.

The **vms\_iowait** function is a special case of the **vms\_iowaitf** function with the **V\_WAITERR** flag set.

## Flags

Name	Purpose
<b>V_WAITERR</b>	Waits until the completion of all I/O unless an error occurs.
<b>V_WAITALL</b>	Waits until the completion of all I/O regardless of any occurrence of I/O errors.

## Execution Environment

The **vms\_iowait** and **vms\_iowaitf** kernel services can be called from the process environment only.

They can only be used by file segments.

## Return Values

Item	Description
0	Indicates that the page-out operations completed.
EIO	Indicates that an error occurred while performing I/O.

### Related reference:

“**vm\_invalidatep** Kernel Service” on page 544

### Related information:

Memory Kernel Services

Understanding Virtual Memory Manager Interfaces

## **vm\_uiomove** Kernel Service

### Purpose

Moves data between a virtual memory object and a buffer specified in the **uio** structure.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
#include <sys/uio.h>
```

```
int vm_uiomove (vmid, limit, rw, uio)  
vmid_t vmid;  
int limit;  
enum uio_rw rw;  
struct uio *uio;
```

## Parameters

Item	Description
<i>vmid</i>	Specifies the virtual memory object ID.
<i>limit</i>	Indicates the limit on the transfer length. If this parameter is negative or 0, no bytes are transferred.
<i>rw</i>	Specifies a read/write flag that gives the direction of the move. The possible values for this parameter ( <b>UIO_READ</b> , <b>UIO_WRITE</b> ) are defined in the <code>/usr/include/sys/uio.h</code> file.
<i>uio</i>	Points to the <b>uio</b> structure.

## Description

The **vm\_uiomove** kernel service moves data between a virtual memory object and the buffer specified in a **uio** structure.

This service determines the virtual addressing required for the data movement according to the offset in the object.

The **vm\_uiomove** kernel service is similar to the **uiomove** kernel service, but the address for the trusted buffer is specified by the *vmid* parameter and the *uio\_offset* field of *offset* parameters instead of as a **caddr\_t** address. The offset size is a 64 bit **offset\_t**, which allows file offsets in client segments which are greater than 2 gigabytes. **vm\_uiomove** must be used instead of **vm\_move** if the client filesystem supports files which are greater than 2 gigabytes.

**Note:** The **vm\_uiomove** kernel service does not support use of cross-memory descriptors.

I/O errors for paging space and a lack of paging space are reported as signals.

## Execution Environment

The **vm\_uiomove** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.
EFAULT	Indicates a bad address.
ENOMEM	Indicates insufficient memory.
ENOSPC	Indicates insufficient disk space.
EIO	Indicates an I/O error.

Other file system-specific **errno** global variables are returned by the virtual file system involved in the move function.

### Related reference:

“uiomove Kernel Service” on page 517

### Related information:

Memory Kernel Services

Understanding Virtual Memory Manager Interfaces



## vm\_umount Kernel Service

### Purpose

Removes a file system from the paging device table.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
int vm_umount ( type, devid)
int type;
dev_t devid ();
```

### Parameters

Item	Description
<i>type</i>	Specifies the type of device. You can specify multiple values. But the <i>type</i> parameter must have a value of <b>D_REMOTE</b> as one of its values. You can also specify the following optional value:  <b>D_NOWAIT</b> Indicates that if I/O discovered during a prior <b>vm_setdevid</b> call has not yet completed, the paging device table entry will be removed, asynchronously, at a future point in time when all such I/O to it has completed. This particular <b>vm_umount</b> kernel service call will return without waiting for the I/O to complete. Any <b>buf</b> structures associated with this paging device entry remain allocated until the paging device entry is finally removed.
<i>devid</i>	Points to the strategy routine.

### Description

The **vm\_umount** kernel service waits for all I/O for the device scheduled by the pager to finish. This service then frees the entry in the paging device table. The associated **buf** structures are also freed.

### Execution Environment

The **vm\_umount** kernel service can be called from the process environment only.

### Return Values

Item	Description
0	Indicates successful completion.
EINVAL	Indicates that a file system with the strategy routine designated by the <i>devid</i> parameter is not in the paging device table.

### Related reference:

“vm\_mount Kernel Service” on page 549

“vm\_setdevid Kernel Service” on page 568

### Related information:

Memory Kernel Services

Understanding Virtual Memory Manager Interfaces

## vm\_write Kernel Service

### Purpose

Initiates page-out for a page range in the address space.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
int vm_write (vaddr, nbytes, force)
int vaddr;
int nbytes;
int force;
```

## Description

The **vm\_write** kernel service initiates page-out for pages that intersect the address range (*vaddr*, *vaddr* + *nbytes*).

If the *force* parameter is nonzero, modified pages are written to disk regardless of how recently they have been written.

Page-out is initiated for each modified page. An unchanged page is left in memory with its reference bit set to 0. This makes the unchanged page a candidate for the page replacement algorithm.

The caller must have write access to the specified pages.

The initiated I/O is asynchronous. The **vms\_iowait** kernel service can be called to wait for I/O completion.

**Note:** The **vm\_write** subroutine is not supported for use on large pages.

## Parameters

Item	Description
<i>vaddr</i>	Specifies the address of the first byte of the page range for which a page-out is desired.
<i>nbytes</i>	Specifies the number of bytes starting at the byte specified by the <i>vaddr</i> parameter. This parameter must be nonnegative. All of the bytes must be in the same virtual memory object.
<i>force</i>	Specifies a flag indicating that a modified page is to be written regardless of when it was last written.

## Execution Environment

The **vm\_write** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful completion.
EINVAL	Indicates one of these four errors: <ul style="list-style-type: none"><li>• A region is not defined.</li><li>• A region is an I/O region.</li><li>• The length specified by the <i>nbytes</i> parameter is negative.</li><li>• The address range crosses a virtual memory object boundary.</li></ul>
EACCES	Indicates that access does not permit writing.
EIO	Indicates a permanent I/O error.

## Related reference:

“vm\_writew Kernel Service” on page 579

“vms\_iowait, vms\_iowaitf Kernel Services” on page 574

## Related information:

Memory Kernel Services

Understanding Virtual Memory Manager Interfaces

## vm\_wri<sup>te</sup>p Kernel Service

### Purpose

Initiates page-out for a page range in a virtual memory object.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
int vm_writep ( vmid, pfirst, npages)
vmid_t vmid;
int pfirst;
int npages;
```

### Description

The **vm\_wri<sup>te</sup>p** kernel service initiates page-out for the specified page range in the virtual memory object. I/O is initiated for modified pages only. Unchanged pages are left in memory, but their reference bits are set to 0.

The caller can wait for the completion of I/O initiated by this and prior calls by calling the **vms\_iowait** kernel service.

**Note:** The **vm\_wri<sup>te</sup>p** subroutine is not supported for use on large pages.

### Parameters

Item	Description
<i>vmid</i>	Specifies the identifier for the virtual memory object.
<i>pfirst</i>	Specifies the first page number at which page-out is to begin.
<i>npages</i>	Specifies the number of pages for which the page-out operation is to be performed.

### Execution Environment

The **vm\_wri<sup>te</sup>p** kernel service can be called from the process environment only.

### Return Values

Item	Description
0	Indicates successful completion.
EINVAL	Indicates any one of the following errors: <ul style="list-style-type: none"><li>• <i>pfirst</i> = 0 and <i>npages</i> = 0.</li><li>• The virtual memory object ID is not valid.</li><li>• The starting page is negative.</li><li>• The number of pages is negative.</li><li>• The page range exceeds the size of virtual memory object.</li><li>• One or more large pages lie in the target page range.</li></ul>

## Related reference:

“vm\_invalidatep Kernel Service” on page 544

“`vm_write` Kernel Service” on page 577

“`vms_iowait`, `vms_iowaitf` Kernel Services” on page 574

**Related information:**

Memory Kernel Services

Understanding Virtual Memory Manager Interfaces

## **vn\_free Kernel Service**

### **Purpose**

Frees a v-node previously allocated by the `vn_get` kernel service.

### **Syntax**

```
#include <sys/types.h>
#include <sys/errno.h>
```

```
int vn_free ( vp)
struct vnode *vp;
```

### **Parameter**

Item	Description
<i>vp</i>	Points to the v-node to be deallocated.

### **Description**

The `vn_free` kernel service provides a mechanism for deallocating v-node objects used within the virtual file system. The v-node specified by the *vp* parameter is returned to the pool of available v-nodes to be used again.

### **Execution Environment**

The `vn_free` kernel service can be called from the process environment only.

### **Return Values**

The `vn_free` service always returns 0.

**Related reference:**

“`vn_get` Kernel Service”

**Related information:**

Virtual File System Overview

Virtual File System (VFS) Kernel Services

## **vn\_get Kernel Service**

### **Purpose**

Allocates a virtual node.

### **Syntax**

```
#include <sys/types.h>
#include <sys/errno.h>
```

```

int vn_get ( vfsp, gnp, vpp)
struct vfs *vfsp;
struct gnode *gnp;
struct vnode **vpp;

```

## Parameters

Item	Description
<i>vfsp</i>	Points to a <b>vfs</b> structure describing the virtual file system that is to contain the v-node. Any returned v-node belongs to this virtual file system.
<i>gnp</i>	Points to the g-node for the object. This pointer is stored in the returned v-node. The new v-node is added to the list of v-nodes in the g-node.
<i>vpp</i>	Points to the place in which to return the v-node pointer. This is set by the <b>vn_get</b> kernel service to point to the newly allocated v-node.

## Description

The **vn\_get** kernel service provides a mechanism for allocating v-node objects for use within the virtual file system environment. A v-node is first allocated from an effectively infinite pool of available v-nodes.

Upon successful return from the **vn\_get** kernel service, the pointer to the v-node pointer provided (specified by the *vpp* parameter) has been set to the address of the newly allocated v-node.

The fields in this v-node have been initialized as follows:

Field	Initial Value
<i>v_count</i>	Set to 1.
<i>v_vfsp</i>	Set to the value in the <i>vfsp</i> parameter.
<i>v_gnode</i>	Set to the value in the <i>gnp</i> parameter.
<i>v_next</i>	Set to list of others v-nodes with the same g-node.

All other fields in the v-node are zeroed.

## Execution Environment

The **vn\_get** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates successful completion.
ENOMEM	Indicates that the <b>vn_get</b> kernel service could not allocate memory for the v-node. (This is a highly unlikely occurrence.)

### Related reference:

“vn\_free Kernel Service” on page 580

### Related information:

Virtual File System Overview

Virtual File System (VFS) Kernel Services

## vsx\_disable Kernel Service

### Purpose

Communicates the status of the vector and the vector-scalar registers to the hypervisor.

## Syntax

```
#include <sys/machine.h>
void vsx_disable (old)
char old;
```

## Parameters

### old

Specifies the value returned by the `vsx_enable` kernel service.

## Description

The `vsx_disable` kernel service communicates to the hypervisor that the vector and the vector-scalar registers are no longer in use.

## Execution Environment

The `vsx_disable` kernel service can be called from the process environment or the interrupt environment. The `vsx_disable` kernel service must be called while all interrupts from within the INTMAX critical section are disabled.

## Return Values

The `vsx_disable` kernel service has no return values.

## vsx\_enable Kernel Service

### Purpose

Communicates the status of the vector and the vector-scalar registers to the hypervisor.

## Syntax

```
#include <sys/machine.h>
char vsx_enable ()
```

## Description

The `vsx_enable` kernel service communicates to the hypervisor that the vector and the vector-scalar registers are in use.

## Execution Environment

The `vsx_enable` kernel service can be called from the process environment or the interrupt environment. The `vsx_enable` kernel service must be called while all the interrupts from within the INTMAX critical section are disabled.

## Return Values

The `vsx_enable` kernel service returns the current setting.

## W

The following kernel services begin with the with the letter w.

## waitcfree Kernel Service

### Purpose

Checks the availability of a free character buffer.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/cblock.h>
#include <sys/sleep.h>
int waitcfree ( )
```

### Description

The `waitcfree` kernel service is used to wait for a buffer which was allocated by a previous call to the `pincf` kernel service. If one is not available, the `waitcfree` kernel service waits until either a character buffer becomes available or a signal is received.

The `waitcfree` kernel service has no parameters.

### Execution Environment

The `waitcfree` kernel service can be called from the process environment only.

### Return Values

Item	Description
EVENT_SUCC	Indicates a successful operation.
EVENT_SIG	Indicates that the wait was terminated by a signal.

### Related reference:

“`pincf` Kernel Service” on page 410

“`putcf` Kernel Service” on page 426

### Related information:

I/O Kernel Services

## waitq Kernel Service

### Purpose

Waits for a queue element to be placed on a device queue.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/deviceq.h>

struct req_qe *waitq ( queue_id)
cba_id queue_id;
```

### Parameter

Item	Description
<i>queue_id</i>	Specifies the device queue identifier.

## Description

The **waitq** kernel service is not part of the base kernel but is provided by the device queue management kernel extension. This queue management kernel extension must be loaded into the kernel before loading any kernel extensions referencing these services.

The **waitq** kernel service waits for a queue element to be placed on the device queue specified by the *queue\_id* parameter. This service performs these two actions:

- Waits on the event mask associated with the device queue.
- Calls the **readq** kernel service to make the most favored queue element the active one.

Processes can only use the **waitq** kernel service to wait for a single device queue. Use the **et\_wait** service to wait on the occurrence of more than one event, such as multiple device queues.

The **waitq** kernel service uses the **EVENT\_SHORT** form of the **et\_wait** kernel service. Therefore, a signal does not terminate the wait. Use the **et\_wait** kernel service if you want a signal to terminate the wait.

The **readq** kernel service can be used to read the active queue element from a queue. It does not wait for a queue element if there are none in the queue.

**Attention:** The server must not alter any fields in the queue element or the system may halt.

## Execution Environment

The **waitq** kernel service can be called from the process environment only.

## Return Values

The **waitq** service returns the address of the active queue element in the device queue.

**Related reference:**

“et\_wait Kernel Service” on page 120

## WPAR\_CKPT\_QUERY (Checkpoint Query) Device Driver ioctl Operation Purpose

Queries a device driver about its checkpoint capabilities.

## Syntax

```
#include <sys/ioctl.h>
```

```
int ioctl ( FileDescriptor, WPAR_CKPT_QUERY, Arg )
```

```
int FileDescriptor;
```

```
wpar_ckpt_resp_t * Arg;
```

## Parameters



Item	Description
<i>FileDescriptor</i>	Open file descriptor that refers to the device being queried for the checkpoint capability.
WPAR_CKPT_QUERY	The command that requests information on the device checkpoint capability.
<i>Arg</i>	Pointer to a <b>wpar_ckpt_resp_t</b> structure which will contain a device driver response on the checkpoint capability upon a successful return from the <b>ioctl</b> call.

## Description

The **WPAR\_CKPT\_QUERY** operation allows a caller to ask a device driver connected to the **ioctl** input file descriptor if it supports checkpoint and restart operations. If a device driver supports checkpoint and restart operations, the returned answer can describe what operations are required to accomplish a checkpoint and restart.

If the device is not checkpoint and restart capable, checkpoint-aware devices fail this **ioctl** request with the **ENOSYS** error. Non-checkpoint-aware devices fail this **ioctl** request as an unknown **ioctl**. If the device is checkpoint and restart capable, checkpoint-aware devices return success.

The *arg* parameter to a **WPAR\_CKPT\_QUERY** **ioctl** request allows the caller to receive specific information regarding how the device supports checkpoint and restart if it is capable. The caller of a **WPAR\_CKPT\_QUERY** **ioctl** request must supply a pointer to a structure of the **wpar\_ckpt\_resp\_t** type in the *arg* parameter.

### wpar\_ckpt\_resp\_t structure

The **wpar\_ckpt\_resp\_t** structure is supplied as the input to the **WPAR\_CKPT\_QUERY** **ioctl** request.

```
#define WPAR_CKPT_OP_MAX 5
typedef struct wpar_ckpt_resp_t {
    int opcnt;
    wpar_ckpt_op_top [WPAR_CKPT_OP_MAX];
}wpar_ckpt_resp_t;
```

The fields of the **wpar\_ckpt\_resp\_t** structure are as follows:

Item	Description
<b>opcnt</b>	Returned from an <b>WPAR_CKPT_QUERY</b> <b>ioctl</b> request as the number of the <b>wpar_ckpt_op_t</b> sub-structures that contain return information.
<b>wpar_ckpt_op_top</b>	A sub-structure that contains specific information on operation types that must occur on a device for it to save or restore its state correctly.

### wpar\_ckpt\_op\_t structure

The **wpar\_ckpt\_op\_t** structure is a sub-structure of the **wpar\_ckpt\_resp\_t** structure.

```
typedef struct wpar_ckpt_op_t {
    int op;
    int opt; /*extended options of openx*/
}wpar_ckpt_op_t;
```

The fields of the **wpar\_ckpt\_op\_t** structure are as follows:

Item	Description
<code>op</code>	Returned from a <code>WPAR_CKPT_QUERY</code> ioctl request. Defined as a set of one or more operations that must be performed to successfully checkpoint and restart the device.
<code>opt</code>	Options to supply to the <code>openx</code> function if the device is to be re-opened on the arrival server through the <code>openx</code> function.

## wpar\_ckpt\_op\_t op field

Item	Description
<code>WPAR_CKPT_OP_NULL</code>	Device requires no special handling for checkpoint and restart operations.
<code>WPAR_CKPT_OP_REOPEN</code>	Device needs to be re-opened through the <code>open</code> function with the access modes applicable at checkpoint time.
<code>WPAR_CKPT_OP_OPENX</code>	Device needs to be re-opened with the <code>openx</code> function. The <code>opt</code> field denotes the desired extension argument to the <code>openx</code> function.

## Return Values

Upon successful completion, this operation returns a value of 0. Otherwise, it returns a value of -1 and the `errno` global variable is set to one of the following values:

Item	Description
<code>ENOSYS</code>	Device cannot participate in checkpoint and restart operations.
<code>EINVAL</code>	Device does not accept the <code>WPAR_CKPT_QUERY</code> operation.

### Related reference:

“`kwpar_checkpoint_status` Kernel Service” on page 314

## w\_clear Kernel Service

### Purpose

Removes a watchdog timer from the list of watchdog timers known to the kernel.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/watchdog.h>
```

```
int w_clear ( w )
struct watchdog *w;
```

### Parameter

Item	Description
<code>w</code>	Specifies the watchdog timer structure.

### Description

The watchdog timer services, including the `w_clear` kernel service, are typically used to verify that an I/O operation completes in a reasonable time.

When the `w_clear` kernel service removes the watchdog timer, the `w->count` watchdog count is no longer decremented. In addition, the `w->func` watchdog timer function is no longer called.

In a uniprocessor environment, the call always succeeds. This is untrue in a multiprocessor environment, where the call will fail if the watchdog timer is being handled by another processor. Therefore, the function now has a return value, which is set to 0 if successful, or -1 otherwise. Funnelled device drivers

do not need to check the return value since they run in a logical uniprocessor environment. Multiprocessor-safe and multiprocessor-efficient device drivers need to check the return value in a loop. In addition, if a driver uses locking, it must release and reacquire its lock within this loop, as shown below:

```
while (w_clear(&watchdog))
    release_then_reacquire_dd_lock;
    /* null statement if locks not used */
```

**Note:** The `w_clear` kernel service clears any attributes that were previously set by using the `w_setattr()` kernel service.

## Execution Environment

The `w_clear` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates that the watchdog timer was successfully removed.
-1	Indicates that the watchdog timer could not be removed.

### Related reference:

“w\_init Kernel Service”

“w\_setattr Kernel Service” on page 588

### Related information:

Timer and Time-of-Day Kernel Services

## w\_init Kernel Service

### Purpose

Registers a watchdog timer with the kernel.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/watchdog.h>
```

```
int w_init ( w )
struct watchdog *w;
```

### Parameter

Item	Description
<i>w</i>	Specifies the watchdog timer structure.

### Description

The `watchdog` structure must be initialized prior to calling the `w_init` kernel service as follows:

- Set the next and prev fields to NULL.
- Set the func and restart fields to the appropriate values.
- Set the count field to 0.

**Attention:** The watchdog structure must be pinned when the `w_init` service is called. It must remain pinned until after the call to the `w_clear` service. During this time, the watchdog structure must not be altered except by the watchdog services.

The watchdog timer services, including the **w\_init** kernel service, are typically used to verify that an I/O operation completes in a reasonable time. The watchdog timer is initialized to the stopped state and must be started using the **w\_start** service.

In both uniprocessor and multiprocessor environments, the **w\_init** kernel service always succeeds.

The calling parameters for the watchdog timer function are:

```
void func (w)
struct watchdog *w;
```

## Execution Environment

The **w\_init** kernel service can be called from the process environment only.

## Return Values

The **w\_init** kernel service returns 0 for compatibility with previous releases of AIX.

### Related reference:

“w\_clear Kernel Service” on page 586

“w\_setattr Kernel Service”

### Related information:

Timer and Time-of-Day Kernel Services

## w\_setattr Kernel Service

### Purpose

Sets attributes for a watchdog timer.

### Syntax

```
#include <sys/watchdog.h>
#include <sys/kerrno.h>
```

```
kerrno_t w_setattr(struct watchdog *w, char attr)
```

### Parameter

Item	Description
<i>w</i>	Specifies the watchdog timer structure.
<i>attr</i>	A bitmask of attributes to be set. Supported flags are: <b>WD_ATTR_MOVE_OK</b> Allow timer to migrate from one CPU to another.

### Description

The **w\_setattr** kernel service sets attributes for the specified watchdog timer. The **WD\_ATTR\_MOVE\_OK** attribute should be set when the caller does not have a dependency on which processor the timer expiration handler is called. This attribute allows the system to move the timer from one processor to another as needed, to improve the effectiveness of processor folding. When this attribute is set, the associated watchdog timer is moved to another processor when the owning processor is folded.

The **w\_setattr** kernel service must be called after the **w\_init()** kernel service but before the **w\_start()** kernel service. Otherwise, the **w\_setattr** kernel service may fail.

## Execution Environment

The `w_setattr` kernel service can be called from either the process or interrupt environment.

## Return Values

Item	Description
0	The specified attribute was successfully set.
<0	The specified attribute was not set. The failure is indicated with return value set to one of the following values:  EINVAL_W_SETATTR_EYEC: An invalid eye catcher was detected.  EINVAL_W_SETATTR_ATTR: An invalid attribute flag was detected.

## Related reference:

“`w_clear` Kernel Service” on page 586

“`w_start` Kernel Service”

“`w_stop` Kernel Service” on page 590

## `w_start` Kernel Service

### Purpose

Starts a watchdog timer.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/watchdog.h>
```

```
void w_start ( w )
struct watchdog *w;
```

### Parameter

Item	Description
<i>w</i>	Specifies the watchdog timer structure.

### Description

The watchdog timers, including the `w_start` kernel service, are typically used to verify that an I/O operation completes in a reasonable time. The `w_start` and `w_stop` kernel services are designed to allow the timer to be started and stopped efficiently. The kernel decrements the `w->count` watchdog count every second. The kernel calls the `w->func` watchdog timer function when the `w->count` watchdog count reaches 0. A watchdog timer is ignored when the `w->count` watchdog count is less than or equal to 0.

The `w_start` kernel service sets the `w->count` watchdog count to a value of `w->restart`.

**Attention:** The watchdog structure must be pinned when the `w_start` kernel service is called. It must remain pinned until after the call to the `w_clear` kernel service. During this time, the watchdog structure must not be altered except by the watchdog services.

## Execution Environment

The `w_start` kernel service can be called from the process and interrupt environments.

## Return Values

The `w_start` kernel service has no return values.

### Related reference:

“`w_stop` Kernel Service”

“`w_setattr` Kernel Service” on page 588

### Related information:

Timer and Time-of-Day Kernel Services

## `w_stop` Kernel Service

### Purpose

Stops a watchdog timer.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/watchdog.h>
```

```
void w_stop ( w)
struct watchdog *w;
```

### Parameter

Item	Description
<i>w</i>	Specifies the watchdog timer structure.

### Description

The watchdog timer services, including the `w_stop` kernel service, are typically used to verify that an I/O operation completes in a reasonable time. The `w_start` and `w_stop` kernel services are designed to allow the timer to be started and stopped efficiently. The kernel decrements the `w->count` watchdog count every second. The kernel calls the `w->func` watchdog timer function when the `w->count` watchdog count reaches 0. A watchdog timer is ignored when `w->count` is less than or equal to 0.

**Attention:** The watchdog structure must be pinned when the `w_stop` kernel service is called. It must remain pinned until after the call to the `w_clear` kernel service. During this time, the watchdog structure must not be altered except by the watchdog services.

### Execution Environment

The `w_stop` kernel service can be called from the process and interrupt environments.

### Return Values

The `w_stop` kernel service has no return values.

### Related reference:

“`w_start` Kernel Service” on page 589

“`w_setattr` Kernel Service” on page 588

### Related information:

Timer and Time-of-Day Kernel Services

## X

The following kernel services begin with the with the letter x.

### **xfidToName() Kernel Service**

#### **Purpose**

Finds the full path name of the file corresponding to an `xfid_t` structure.

#### **Syntax**

```
#include <sys/xfops.h>
```

```
int    xfidToName(struct xfid *xfp,  
                  void *nrp,  
                  char *pathname,  
                  unsigned int pbufLen,  
                  long flags);
```

#### **Description**

The `xfidToName()` kernel service finds a name for an `xfid` value.

#### **Parameters**

##### **xfp**

Pointer to the `xfid` value for which a name is needed.

##### **nrp**

Name resolution structure that is passed to the validation routine.

##### **pathname**

Pointer to buffer where the file name will be stored.

##### **pbufLen**

Size of path name buffer. A size of `MAXPATHLEN` is sufficient to hold any path name.

##### **flags**

Operation modifiers. This parameter must be set to zero.

#### **Return values**

**0** Indicates success. The path name for the `xfid` value is returned.

##### **ENOENT**

Name not found.

##### **EPERM**

No permission for lookup.

##### **EINVAL**

Invalid parameter is specified.

**E2BIG** Path is larger than `pbufLen` bytes.

### **xlata\_create Kernel Service**

#### **Purpose**

Creates pretranslation data structures.

## Syntax

```
int xlate_create (dp, baddr, count, flags)
struct xmem*dp;
caddr_t baddr;
int count;
uint flags;
```

## Description

The `xlate_create` kernel service creates pretranslation data structures capable of pretranslating all pages of the virtual buffer indicated by the `baddr` parameter for length of `count` into a list of physical page numbers, appended to the cross memory descriptor pointed to by `dp`.

If the `XLATE_ALLOC` flag is set, only the data structures are created and no pretranslation is done. If the flag is not set, in addition to the data structures being created, each page of the buffer is translated and the access permissions verified, requiring read-write access to each page. The `XLATE_ALLOC` flag is useful when the buffer will be pinned and utilized later, through the `xlate_pin` and `xlate_unpin` kernel services.

The `XLATE_SPARSE` flag can be used to indicate that only selected portions of a pretranslated region may be valid (pinned and pretranslated) at any given time. The `XLATE_SPARSE` flag can be used in conjunction with the `XLATE_ALLOC` flag to preallocate the pretranslation data structures for an address region that will be dynamically managed.

The `xlate_create` kernel service is primarily for use when memory buffers will be reused for I/O. The use of this service to create a pretranslation for the memory buffer avoids page translation and access checking overhead for all future DMAs involving the memory buffer until the `xlate_remove` kernel service is called.

## Parameters

Item	Description
<i>dp</i>	Points to the cross memory descriptor.
<i>baddr</i>	Points to the virtual buffer.
<i>count</i>	Specifies the length of the virtual buffer.
<i>flags</i>	Specifies the operation. Valid values are as follows:  <b>XLATE_PERSISTENT</b> Indicates that the pretranslation data structures should be persistent across calls to pretranslation services.  <b>XLATE_ALLOC</b> Indicates that the pretranslation data structures should be allocated only, and no translation should be performed.  <b>XLATE_SPARSE</b> Indicates that the pretranslation information will be sparse, allowing for the coexistence of valid (active) pretranslation regions and invalid (inactive) pretranslation regions.

## Return Values



Item	Description
ENOMEM	Unable to allocate memory
XMEM_FAIL	No physical translation, or No Access to a Page
XMEM_SUCC	Successful pretranslation created

## Execution Environment

The `xlate_create` kernel service can only be called from the process environment. The entire buffer must be pinned (unless the `XLATE_ALLOC` flag is set), and the cross memory descriptor valid.

### Related reference:

“`xlate_remove` Kernel Service” on page 594

“`xlate_pin` Kernel Service”

“`xlate_unpin` Kernel Service” on page 595

## `xlate_pin` Kernel Service

### Purpose

Pins all pages of a virtual buffer.

### Syntax

```
int xlate_pin (dp, baddr, count, rw)
struct xmem *dp;
caddr_t baddr;
int count;
int rw;
```

### Description

The `xlate_pin` kernel service pins all pages of the virtual buffer indicated by the `baddr` parameter for length of `count` and also appends pretranslation information to the cross memory descriptor pointed to by the `dp` parameter.

The `xlate_pin` kernel service results in a short-term pin, which will support `mmap` and `shmat` allocated memory buffers.

In addition to pinning and translating each page, the access permissions to the page are verified according to the desired access (as specified by the `rw` parameter). For a setting of `B_READ`, write access to the page must be allowed. For a setting of `B_WRITE`, only read access to the page must be allowed.

The caller can preallocate pretranslation data structures and append them to the cross memory descriptor prior to the call (through a call to the `xlate_create` kernel service), or have this service allocate the necessary data structures. If the cross memory descriptor is already of type `XMEM_XLATE`, it is assumed that the data structures are already allocated. If callers want to have the pretranslation data structures persist across the subsequent `xlate_unpin` call, they should also set the `XLATE_PERSISTENT` flag on the call to the `xlate_create` kernel service.

### Parameters

Item	Description
<i>dp</i>	Points to the cross memory descriptor.
<i>baddr</i>	Points to the virtual buffer.
<i>count</i>	Specifies the length of the virtual buffer.
<i>rw</i>	Specifies the access permissions for each page.

## Return Values

If successful, the `xlate_pin` kernel service returns 0. If unsuccessful, one of the following is returned:

Item	Description
EINVAL	Invalid cross memory descriptor or parameters.
ENOMEM	Unable to allocate memory.
ENOSPC	Out of Paging Resources.
XMEM_FAIL	Page Access violation.

## Execution Environment

The `xlate_pin` kernel service is only callable from the process environment, and the cross memory descriptor must be valid.

### Related reference:

“`xm_det` Kernel Service” on page 596

“`xm_mapin` Kernel Service” on page 596

“`xlate_unpin` Kernel Service” on page 595

## `xlate_remove` Kernel Service

### Purpose

Removes physical translation information from an `xmem` descriptor from a prior `xlate_create` call.

### Syntax

```
caddr_t xlate_remove (dp)
struct xmem *dp;
```

### Description

See the `xlate_create` kernel service.

### Parameters

Item	Description
<i>dp</i>	Points to the cross memory descriptor.

## Return Values

Item	Description
XMEM_FAIL	No pretranslation information present in the xmem descriptor.
XMEM_SUCC	Pretranslation successfully removed.

## Execution Environment

The `xlate_remove` kernel service can only be called from the process environment.

### Related reference:

“xm\_det Kernel Service” on page 596

“xlate\_pin Kernel Service” on page 593

“xlate\_unpin Kernel Service”

## xlate\_unpin Kernel Service Purpose

Unpins all pages of a virtual buffer.

### Syntax

```
int xlate_unpin (dp, baddr, count)
struct xmem *dp;
caddr_t baddr;
int count;
```

### Description

The `xlate_unpin` kernel service unpins pages from a prior call to the `xlate_pin` kernel service based on the `baddr` and `count` parameters. It does this by utilizing the pretranslated real page numbers appended to the cross memory descriptor pointed to by `dp`.

If the `XLATE_PERSISTENT` flag is not set in the `prexflags` flag word of the pretranslation data structure, the pretranslation data structures are also freed.

### Parameters

Item	Description
<code>dp</code>	Points to the cross memory descriptor.
<code>baddr</code>	Points to the virtual buffer.
<code>count</code>	Specifies the length of the virtual buffer.

### Return Values

If successful, the `xlate_unpin` kernel service returns 0. If unsuccessful, one of the following is returned:

Item	Description
EINVAL	Invalid cross memory descriptor or parameters.
ENOSPC	Unable to allocate paging space (case of <code>mmap</code> segment).
ENOSPC	Out of Paging Resources.
XMEM_FAIL	Page Access violation.

### Related reference:

“xm\_det Kernel Service” on page 596

“xm\_mapin Kernel Service” on page 596

“xlate\_pin Kernel Service” on page 593

## xm\_det Kernel Service

### Purpose

Releases the addressability to the address space described by an xmem descriptor.

### Syntax

```
void xm_det (baddr, dp)
caddr_t baddr;
struct xmem *dp;
```

### Description

See the **xm\_mapin** Kernel Service for more information.

### Parameters

Item	Description
<i>baddr</i>	Specifies the effective address previously returned from the <b>xm_mapin</b> kernel service.
<i>dp</i>	Cross memory descriptor that describes the above memory object.

### Related reference:

“xlate\_create Kernel Service” on page 591

“xlate\_remove Kernel Service” on page 594

“xm\_mapin Kernel Service”

## xm\_mapin Kernel Service

### Purpose

Sets up addressability in the current process context.

### Syntax

```
#include <sys/adspace.h>int xm_mapin (dp, baddr, count, eaddr)
struct xmem *dp; caddr_t baddr;
size_t count;
caddr_t *eaddr;
```

### Description

The **xm\_mapin** kernel service sets up addressability in the current process context to the address space indicated by the cross memory descriptor pointed to by the *dp* parameter for the addresses [*baddr*, *baddr* + *count* - 1].

This service is created specifically for Client File Systems, or others who need to setup addressability to an address space defined by an xmem descriptor.

If the requested mapping spans a segment boundary, no mapping will be performed, and a return code of **EAGAIN** is returned to indicate that individual calls to the **xm\_mapin** kernel service are necessary to map the portions of the buffer in each segment. The **xm\_mapin** kernel service must be called again with the original *baddr* and a *count* indicating the number of bytes to the next segment. (The number of bytes to the next segment boundary can be obtained using the **xm\_maxmap** kernel service.) This will provide an effective address to use for accessing this portion of the buffer. Then, iteratively, **xm\_mapin** must be called with the segment boundary address (previous *baddr* + *count*), and a new *count* indicating the remainder of the buffer or the next segment boundary, whichever is smaller. This will provide another effective address to use for accessing the next portion of the buffer.

Each address set up by the **xm\_mapin** kernel service must be undone with the **xm\_det** kernel service when it is no longer needed because the **xm\_mapin** kernel service currently uses the **vm\_att** kernel service.

## Parameters

Item	Description
<i>dp</i>	Points to the cross memory descriptor.
<i>baddr</i>	Points to the virtual buffer.
<i>count</i>	Specifies the length of the virtual buffer to map.
<i>eaddr</i>	Points to where the effective address to access the data buffer is returned.

## Return Values

Item	Description
0	Successful. (Reference Parameter <i>eaddr</i> contains the address to use)
XMEM_FAIL	Invalid cross memory descriptor.
EAGAIN	Segment boundary crossing encountered. Caller should make separate <b>xm_mapin</b> calls to map each segments worth.

## Execution Environment

The **xm\_mapin** kernel service can be called from the process or interrupt environments.

### Related reference:

- “xm\_det Kernel Service” on page 596
- “xlate\_remove Kernel Service” on page 594
- “xlate\_pin Kernel Service” on page 593
- “xm\_maxmap Kernel Service”

## xm\_maxmap Kernel Service

### Purpose

Determines the maximum permissible count value for a subsequent call to **xm\_mapin**.

### Syntax

```
#include <sys/adspace.h>
```

```
int xm_maxmap (dp, uaddr, len)
  struct xmem *dp;
  void *uaddr;
  size_t *len;
```

### Parameters

Item	Description
<i>dp</i>	Points to the cross memory descriptor.
<i>uaddr</i>	Points to the virtual buffer.
<i>len</i>	Points to where the maximum permissible count is returned.

### Description

The **xm\_maxmap** kernel service determines the maximum permissible count value (in bytes) for a subsequent **xm\_mapin** call. The value is determined based on the input cross-memory descriptor *dp* and the starting address *uaddr*, and it is returned in the *len* parameter. There is no guarantee that **xm\_mapin**

will succeed; however, it is guaranteed that *uaddr + \*len - 1* is in the same segment as *uaddr*, and therefore *xm\_mapin* will not return **EAGAIN**.

## Execution Environment

The *xm\_maxmap* interface can be called from the process or interrupt environment.

## Return Values

Item	Description
<b>XMEM_SUCC</b>	Successful (Reference parameter <i>len</i> contains the maximum permissible value for a subsequent <i>xm_mapin</i> call)
<b>XMEM_FAIL</b>	Invalid cross memory descriptor.
<b>EAGAIN</b>	Segment boundary crossing encountered. Caller should make separate <i>xm_mapin</i> calls to map each segment's worth.

## Related reference:

"*xm\_mapin* Kernel Service" on page 596

## xmalloc Kernel Service

### Purpose

Allocates memory.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/malloc.h>
```

```
caddr_t xmalloc ( size, align, heap)
int size;
int align;
caddr_t heap;
```

### Parameters

Item	Description
<i>size</i>	Specifies the number of bytes to allocate.
<i>align</i>	Specifies the alignment characteristics for the allocated memory.
<i>heap</i>	Specifies the address of the heap from which the memory is to be allocated.

### Description

The *xmalloc* kernel service allocates an area of memory out of the heap specified by the *heap* parameter. This area is the number of bytes in length specified by the *size* parameter and is aligned on the byte boundary specified by the *align* parameter. The *align* parameter is actually the log base 2 of the desired address boundary. For example, an *align* value of 4 requests that the allocated area be aligned on a 2<sup>4</sup> (16) byte boundary.

There are multiple heaps provided by the kernel for use by kernel extensions. Two primary kernel heaps are **kernel\_heap** and **pinned\_heap**. Kernel extensions should use the **kernel\_heap** value when allocating memory that is not pinned, and should use the **pinned\_heap** value when allocating memory that should always be pinned or pinned for long periods of time. When allocating from the **pinned\_heap** heap, the *xmalloc* kernel service will pin the memory before a successful return. The **pin** and **unpin** kernel services should be used to pin and unpin memory from the **kernel\_heap** heap when the memory should only be pinned for a limited amount of time. Memory from the **kernel\_heap** heap must be unpinned before freeing it. Memory from the **pinned\_heap** heap should not be unpinned.

The **kernel\_heap** heap points to one of the following heaps: **kernel\_heap\_4K\_64K** and **kernel\_heap\_16M**. The **pinned\_heap** heap points to one of the following heaps: **pinned\_heap\_4K\_64K** and **pinned\_heap\_16M**. Each of the target heaps differ in the size of the pages that back them. **kernel\_heap\_4K\_64K** or **pinned\_heap\_4K\_64K** will be backed by either medium (64 KB) or regular (4 KB) pages, depending on the page size supported by the machine. **kernel\_heap\_16M** or **pinned\_heap\_16M** will return memory backed by large pages if large page heaps are enabled. If large page heaps are not enabled, **kernel\_heap** or **pinned\_heap** will point to the default heap. If the size of the backing pages are not important, use the **kernel\_heap** value and the **pinned\_heap** value. They will point to the heap that you prefer. For more information about large page heap support, see **vmo**.

Kernel extensions can use these services to allocate memory out of the kernel heaps. For example, the **xmalloc (128,3,kernel\_heap)** kernel service allocates a 128-byte double word aligned area out of the kernel heap.

A kernel extension must use the **xmfree** kernel service to free the allocated memory. If it does not, subsequent allocations eventually are unsuccessful.

The **xmalloc** kernel service has two compatibility interfaces: **malloc** and **palloc**.

The following additional interfaces to the **xmalloc** kernel service are provided:

- **malloc** (*size*) is equivalent to **xmalloc** (*size*, 0, **kernel\_heap**).
- **palloc** (*size*, *align*) is equivalent to **xmalloc** (*size*, *align*, **kernel\_heap**).

## Execution Environment

The **xmalloc** kernel service can be called from the process environment only.

## Return Values

Upon successful completion, the **xmalloc** kernel service returns the address of the allocated area. A null pointer is returned under the following circumstances:

- The requested memory cannot be allocated.
- The heap has not been initialized for memory allocation.

### Related reference:

“xmfree Kernel Service” on page 611

### Related information:

Memory Kernel Services

## xmattach Kernel Service

### Purpose

Attaches to a user buffer for cross-memory operations.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/xmem.h>
```

```
int xmattach (addr, count, dp, segflag)
char * addr;
int count;
struct xmem * dp;
int segflag;
```

## Parameters

Item	Description
<i>addr</i>	Specifies the address of the user buffer to be accessed in a cross-memory operation.
<i>count</i>	Indicates the size of the user buffer to be accessed in a cross-memory operation.
<i>dp</i>	Specifies a cross-memory descriptor. The <i>dp-&gt;aspace_id</i> variable must be set to a value of <b>XMEM_INVALID</b> .
<i>segflag</i>	Specifies a segment flag. This flag is used to determine the address space of the memory that the cross-memory descriptor applies to, as well as for other purposes. The valid values for this flag can be found in the <code>/usr/include/xmem.h</code> file.

## Description

The **xmattach** kernel service prepares the user buffer so that a device driver can access it without executing under the process that requested the I/O operation. A device top-half routine calls the **xmattach** kernel service. The **xmattach** kernel service allows a kernel process or device bottom-half routine to access the user buffer with the **xmemin** or **xmemout** kernel services. The device driver must use the **xmdetach** kernel service to inform the kernel when it has finished accessing the user buffer.

The kernel remembers which segments are attached for cross-memory operations. Resources associated with these segments cannot be freed until all cross-memory descriptors have been detached. "Cross Memory Kernel Services" in Memory Kernel Services in *Kernel Extensions and Device Support Programming Concepts* describes how the cross-memory kernel services use cross-memory descriptors.

**Note:** When the **xmattach** kernel service remaps user memory containing the cross-memory buffer, the effects are machine-dependent. Also, cross-memory descriptors are not inherited by a child process.

Storage-key protection can be enforced on memory regions described by a cross-memory descriptor. The enforcement is done during normal access checking performed by cross-memory services, such as the **xmmdma** kernel service. A kernel keyset can be contained in the cross-memory descriptor to limit memory accessibility. When a keyset is associated with a cross-memory descriptor, access to the memory region is limited by that keyset. A keyset is required because a cross-memory descriptor can describe a virtual memory region with multiple keys assigned to the pages it contains. Normally, a keyset describes the accessibility of the context that the attach was initiated for. For example, a cross-memory attached to user-space contains a description of the user-mode accessibility (keyset). Adding keysets to kernel cross-memory descriptors can also enhance system RAS, since they limit kernel access by the cross-memory descriptor. Typically it is limited to that of the **xmattach** caller or to specific key(s), to catch cases where a cross-memory descriptor is misused.

User-mode storage-keys are always associated with descriptors attached using **USER\_SPACE** or **USER1\_SPACE** segflag. These flags were always required to attach to the user address space, so no explicit update is required to enable storage-key protection on user memory attaches. Once attached, existing kernel services that require cross-memory descriptors enforce the user keyset saved at attach time when performing memory accesses or checking user accessibility.

For kernel memory, a keyset is not used to restrict regions attached with **SYS\_ADSPACE**. Attaching a region with **SYS\_ADSPACE\_ASSIGN\_KEYSET** associates the caller's keyset with the cross-memory region.

## Execution Environment

The **xmattach** kernel service can be called from the process environment only.

## Return Values



Item	Description
XMEM_SUCC	Indicates a successful operation.
XMEM_FAIL	Indicates one of the following errors: <ul style="list-style-type: none"> <li>• The buffer size indicated by the <i>count</i> parameter is less than or equal to 0.</li> <li>• The cross-memory descriptor is in use (<i>dp-&gt;aspace_id</i> != XMEM_INVALID).</li> <li>• The area of memory indicated by the <i>addr</i> and <i>count</i> parameters is not defined.</li> </ul>

#### Related reference:

“uphysio Kernel Service” on page 524

“xmdetach Kernel Service”

“xmgethkeyset Kernel Service” on page 612

## xmdetach Kernel Service

### Purpose

Detaches from a user buffer used for cross-memory operations.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/xmem.h>
```

```
int xmdetach ( dp )
struct xmem *dp;
```

### Parameter

Item	Description
<i>dp</i>	Points to a cross-memory descriptor initialized by the <b>xmattach</b> kernel service.

### Description

The **xmdetach** kernel service informs the kernel that a user buffer can no longer be accessed. This means that some previous caller, typically a device driver bottom half or a kernel process, is no longer permitted to do cross-memory operations on this buffer. Subsequent calls to either the **xmemin** or **xmemout** kernel service using this cross-memory descriptor result in an error return. The cross-memory descriptor is set to *dp->aspace\_id* = XMEM\_INVALID so that the descriptor can be used again. "Cross Memory Kernel Services" in Memory Kernel Services in *Kernel Extensions and Device Support Programming Concepts* describes how the cross-memory kernel services use cross-memory descriptors.

### Execution Environment

The **xmdetach** kernel service can be called from either the process or interrupt environment.

### Return Values

Item	Description
XMEM_SUCC	Indicates successful completion.
XMEM_FAIL	Indicates that the descriptor was not valid or the buffer was not defined.

#### Related reference:

“xmattach Kernel Service” on page 599

“xmemout Kernel Service” on page 609

#### Related information:

Cross Memory Kernel Services

## xmemdma Kernel Service

### Purpose

Prepares a page for direct memory access (DMA) I/O or processes a page after DMA I/O is complete.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/xmem.h>
```

```
int xmemdma ( xp, xaddr, flag)
struct xmem *xp;
caddr_t xaddr;
int flag;
```

### Parameters

Item	Description
<i>xp</i>	Specifies a cross-memory descriptor.
<i>xaddr</i>	Identifies the address specifying the page for transfer.
<i>flag</i>	Specifies whether to prepare a page for DMA I/O or process it after DMA I/O is complete. Possible values are: <ul style="list-style-type: none"> <li><b>XMEM_ACC_CHK</b> Performs access checking on the page. When this flag is set, the page protection attributes are verified.</li> <li><b>XMEM_DR_SAFE</b> Indicates that the use of the real memory address is DLPAR safe.</li> <li><b>XMEM_HIDE</b> Prepares the page for DMA I/O. For cache-inconsistent platforms, this preparation includes hiding the page by making it inaccessible.</li> <li><b>XMEM_UNHIDE</b> Processes the page after DMA I/O. Also, this flag reveals the page and makes it accessible for cache-inconsistent platforms.</li> <li><b>XMEM_WRITE_ONLY</b> Marks the intended transfer as outbound only. This flag is used with <b>XMEM_ACC_CHK</b> to indicate that read-only access to the page is sufficient.</li> </ul>

### Description

The **xmemdma** kernel service operates on the page specified by the *xaddr* parameter in the region specified by the cross-memory descriptor. If the cross-memory descriptor is for the kernel, the *xaddr* parameter specifies a kernel address. Otherwise, the *xaddr* parameter specifies the offset in the region described in the cross-memory descriptor.

The **xmemdma** kernel service is provided for machines that have processor-memory caches, but that do not perform DMA I/O through the cache. Device handlers for Micro Channel DMA devices use the **d\_master** service and **d\_complete** kernel service instead of the **xmemdma** kernel service.

If the *flag* parameter indicates `XMEM_HIDE` (that is, `XMEM_UNHIDE` is not set) and this is the first hide for the page, the `xmemdma` kernel service prepares the page for DMA I/O by flushing the cache and making the page invalid. When the `XMEM_UNHIDE` bit is set and this is the last unhide for the page, the following events take place:

1. The page is made valid.  
If the page is not in pager I/O state:
2. Any processes waiting on the page are readied.
3. The modified bit for the page is set unless the page has a read-only storage key.

The page is made not valid during DMA operations so that it is not addressable with any virtual address. This prevents any process from reading or loading any part of the page into the cache during the DMA operation.

The page specified must be in memory and must be pinned.

If the `XMEM_ACC_CHK` bit is set, then the `xmemdma` kernel service also verifies access permissions to the page. If the page access is read-only, then the `XMEM_WRITE_ONLY` bit must be set in the *flag* parameter.

**Note:**

1. The `xmemdma` kernel service does not hide or reveal the page nor does it perform any cache flushing. The service's primary function is for real-address translation.
2. This service is not supported for large-memory systems with greater than 4GB of physical memory addresses. For such systems, `xmemdma64` should be used.

## Execution Environment

The `xmemdma` kernel service can be called from either the process or interrupt environment.

## Return Values

On successful completion, the `xmemdma` service returns the real address corresponding to the *xaddr* and *xp* parameters.

## Error Codes

The `xmemdma` kernel service returns a value of `XMEM_FAIL` if one of the following are true:

- The descriptor was invalid.
- The page specified by the *xaddr* or *xp* parameter is invalid.
- Access is not allowed to the page.

**Related information:**

Cross Memory Kernel Services

Understanding Direct Memory Access (DMA) Transfer

Dynamic Logical Partitioning

## xmemdma64 Kernel Service

### Purpose

Prepares a page for direct memory access (DMA) I/O or processes a page after DMA I/O is complete.

## Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/xmem.h>

unsigned long long xmemdma64 (
struct xmem *dp,
caddr_t xaddr,>
int flags)
```

## Parameters

Item	Description
<i>dp</i>	Specifies a cross-memory descriptor.
<i>xaddr</i>	Identifies the address that specify the page for transfer.
<i>flags</i>	Specifies whether to prepare a page for DMA I/O or process it after DMA I/O is complete. Possible values are:  <b>XMEM_HIDE</b> Prepares the page for DMA I/O. If cache-inconsistent, then the data cache is flushed, the memory page is hidden, and the real page address is returned. If cache-consistent, then the modified bit is set and the real address of the page is returned.  <b>XMEM_UNHIDE</b> Processes the page after DMA I/O. Also, this flag reveals the page, readies any waiting processes on the page, and sets the modified bit accordingly.  <b>XMEM_ACC_CHK</b> Performs access checking on the page. When this flag is set, the page protection attributes are verified.  <b>XMEM_WRITE_ONLY</b> Marks the intended transfer as outbound only. This flag is used with <b>XMEM_ACC_CHK</b> to indicate that read-only access to the page is sufficient.

## Description

The **xmemdma64** kernel service operates on the page that is specified by the *xaddr* parameter in the region that is specified by the cross-memory descriptor. If the cross-memory descriptor is for the kernel, the *xaddr* parameter specifies a kernel address. Otherwise, the *xaddr* parameter specifies the offset in the region that is described in the cross-memory descriptor.

The **xmemdma64** kernel service is provided for machines that have processor-memory caches, but that do not perform DMA I/O through the cache.

If the *flag* parameter indicates **XMEM\_HIDE** (that is, **XMEM\_UNHIDE** is not set) and it is the first hide for the page, the **xmemdma64** kernel service prepares the page for DMA I/O by flushing the cache and making the page invalid. When the **XMEM\_UNHIDE** bit is set and it is the last unhide for the page, the following events take place:

1. The page is made valid.  
If the page is not in pager I/O state:
2. Any processes that is waiting on the page are readied.
3. The modified bit for the page is set unless the page has a read-only storage key.

The page is made not valid during DMA operations so that it is not addressable with any virtual address. It prevents any process from reading or loading any part of the page into the cache during the DMA operation.

The page that is specified must be in memory and must be pinned.

If the `XMEM_ACC_CHK` bit is set, then the `xmemdma64` kernel service also verifies access permissions to the page. If the page access is read-only, then the `XMEM_WRITE_ONLY` bit must be set in the *flag* parameter.

**Note:** The `xmemdma64` kernel service does not hide or reveal the page, nor does it perform any cache flushing. The service's primary function is for real-address translation.

## Execution Environment

The `xmemdma64` kernel service can be called from either the process or interrupt environment.

## Return Values

On successful completion, the `xmemdma64` service returns the real address corresponding to the *xaddr* and *xp* parameters.

## Error Codes

The `xmemdma64` kernel service returns a value of `XMEM_FAIL` if one of the following are true:

- The descriptor was invalid.
- The page that is specified by the *xaddr* or *xp* parameter is invalid.
- Access is not allowed to the page.

### Related information:

Cross Memory Kernel Services

Understanding Direct Memory Access (DMA) Transfer

## xmempin Kernel Service

### Purpose

Pins the specified address range in user or system memory.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/uio.h>
```

```
int xmempin( base, len, xd)
caddr_t base;
int len;
struct xmem *xd;
```

### Parameters

Item	Description
<i>base</i>	Specifies the address of the first byte to pin.
<i>len</i>	Indicates the number of bytes to pin.
<i>xd</i>	Specifies the cross-memory descriptor.

### Description

The `xmempin` kernel service is used to pin pages backing a specified memory region which is defined in either system or user address space. Pinning a memory region prohibits the pager from stealing pages from the pages backing the pinned memory region. Once a memory region is pinned, accessing that region does not result in a page fault until the region is subsequently unpinned.

The cross-memory descriptor must have been filled in correctly prior to the **xmempin** call (for example, by calling the **xmattach** kernel service).

## Execution Environment

The **xmempin** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates successful completion.
EFAULT	Indicates that the memory region as specified by the <i>base</i> and <i>len</i> parameters is not within the address space specified by the <i>xd</i> parameter.
EINVAL	Indicates that the value of the length parameter is negative or 0. Otherwise, the area of memory beginning at the byte specified by the <i>base</i> parameter and extending for the number of bytes specified by the <i>len</i> parameter is not defined.
ENOMEM	Indicates that the <b>xmempin</b> kernel service is unable to pin the region due to insufficient real memory or because it has exceeded the systemwide pin count.

### Related reference:

“pin Kernel Service” on page 408

“xmemunpin Kernel Service”

### Related information:

Memory Kernel Services

## xmemunpin Kernel Service

### Purpose

Unpins the specified address range in user or system memory.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/uio.h>
```

```
int xmemunpin ( base, len, xd)
caddr_t base;
int len;
struct xmem *xd;
```

### Parameters

Item	Description
<i>base</i>	Specifies the address of the first byte to unpin.
<i>len</i>	Indicates the number of bytes to unpin.
<i>xd</i>	Specifies the cross-memory descriptor.

### Description

The **xmemunpin** kernel service unpins a region of memory. When the pin count is 0, the page is not pinned and can be paged out of real memory. Upon finding an unpinned page, the **xmemunpin** kernel service returns the **EINVAL** error code and leaves any remaining pinned pages still pinned.

The **xmemunpin** service should be used where the address space might be in either user or kernel space.

The cross-memory descriptor must have been filled in correctly prior to the **xmempin** call (for example, by calling the **xmattach** kernel service).

## Execution Environment

The `xmemunpin` kernel service can be called in the process environment when unpinning data that is in either user space or system space. It can be called in the interrupt environment only when unpinning data that is in system space.

## Return Values

Item	Description
0	Indicates successful completion.
EFAULT	Indicates that the memory region as specified by the <i>base</i> and <i>len</i> parameters is not within the address specified by the <i>xd</i> parameter.
EINVAL	Indicates that the value of the length parameter is negative or 0. Otherwise, the area of memory beginning at the byte specified by the <i>base</i> parameter and extending for the number of bytes specified by the <i>len</i> parameter is not defined. If neither cause is responsible, an unpinned page was specified.

### Related reference:

“unpin Kernel Service” on page 520

“xmempin Kernel Service” on page 605

### Related information:

Understanding Execution Environments

## xmemzero Kernel Service

### Purpose

Zeros a buffer described by a cross memory descriptor.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/vmuser.h>
```

```
int xmemzero ( dp, uaddr, count)
struct xmem * dp;
caddr_t uaddr;
long count;
```

### Parameters

Item	Description
<i>dp</i>	The cross memory descriptor.
<i>uaddr</i>	The address in the buffer to begin zeroing.
<i>count</i>	The number of bytes to be zeroed.

### Description

The `xmemzero` kernel service zeros a buffer described by a cross memory descriptor. The page specified must be in memory.

## Execution Environment

The `xmemzero` kernel service can be called from a process or an interrupt environment.

## Return Values

Item	Description
XMEM_SUCC	Indicates the area in the buffer has been zeroed.
XMEM_FAIL	Indicates one of the following errors: <ul style="list-style-type: none"> <li>• The descriptor is marked by XMEM_REMIO.</li> <li>• The descriptor is not marked by XMEM_PROC and XMEM_GLOBAL.</li> <li>• Count &lt; 0.</li> </ul>

### Related information:

Memory Kernel Services

Understanding Virtual Memory Manager Interfaces

## xmemin Kernel Service

### Purpose

Performs a cross-memory move by copying data from the specified address space to kernel global memory.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/xmem.h>
```

```
int xmemin (uaddr, kaddr, count, dp)
caddr_t * uaddr;
caddr_t * kaddr;
int count;
struct xmem * dp;
```

### Parameters

Item	Description
<i>uaddr</i>	Specifies the address in memory specified by a cross-memory descriptor.
<i>kaddr</i>	Specifies the address in kernel memory.
<i>count</i>	Specifies the number of bytes to copy.
<i>dp</i>	Specifies the cross-memory descriptor.

### Description

The **xmemin** kernel service performs a cross-memory move. A cross-memory move occurs when data is moved to or from an address space other than the address space that the program is executing in. The **xmemin** kernel service copies data from the specified address space to kernel global memory.

The **xmemin** kernel service is provided so that kernel processes and interrupt handlers can safely access a buffer within a user process. Calling the **xmattach** kernel service prepares the user buffer for the cross-memory move.

The **xmemin** kernel service differs from the **copyin** and **copyout** kernel services in that it is used to access a user buffer when not executing under the user process. In contrast, the **copyin** and **copyout** kernel services are used only to access a user buffer while executing under the user process.

### Execution Environment

The **xmemin** kernel service can be called from either the process or interrupt environment.



## Return Values

Item	Description
<code>XMEM_SUCC</code>	Indicates successful completion.
<code>XMEM_FAIL</code>	Indicates one of the following errors: <ul style="list-style-type: none"><li>• The user does not have the appropriate access authority for the user buffer.</li><li>• The user buffer is located in an address range that is not valid.</li><li>• The segment containing the user buffer has been deleted.</li><li>• The cross-memory descriptor is not valid.</li><li>• A paging I/O error occurred while the user buffer was being accessed.</li></ul> If the user buffer is not in memory, the <code>xmemin</code> kernel service also returns an <code>XMEM_FAIL</code> error when executing on an interrupt level.

### Related reference:

“`xmattach` Kernel Service” on page 599

“`xmemout` Kernel Service”

### Related information:

Cross Memory Kernel Services

## xmemout Kernel Service

### Purpose

Performs a cross-memory move by copying data from kernel global memory to a specified address space.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/xmem.h>
```

```
int xmemout (kaddr, uaddr, count, dp)
caddr_t * kaddr;
caddr_t * uaddr;
int count;
struct xmem * dp;
```

### Parameters

Item	Description
<i>kaddr</i>	Specifies the address in kernel memory.
<i>uaddr</i>	Specifies the address in memory specified by a cross-memory descriptor.
<i>count</i>	Specifies the number of bytes to copy.
<i>dp</i>	Specifies the cross-memory descriptor.

### Description

The `xmemout` kernel service performs a cross-memory move. A cross-memory move occurs when data is moved to or from an address space other than the address space that the program is executing in. The `xmemout` kernel service copies data from kernel global memory to the specified address space.

The `xmemout` kernel service is provided so that kernel processes and interrupt handlers can safely access a buffer within a user process. Calling the `xmattach` kernel service prepares the user buffer for the cross-memory move.

The `xmemout` kernel service differs from the `copyin` and `copyout` kernel services in that it is used to access a user buffer when not executing under the user process. In contrast, the `copyin` and `copyout`

kernel services are only used to access a user buffer while executing under the user process.

## Execution Environment

The `xmemout` kernel service can be called from either the process or interrupt environment.

## Return Values

Item	Description
<code>XMEM_SUCC</code>	Indicates successful completion.
<code>XMEM_FAIL</code>	Indicates one of the following errors: <ul style="list-style-type: none"><li>• The user does not have the appropriate access authority for the user buffer.</li><li>• The user buffer is located in an address range that is not valid.</li><li>• The segment containing the user buffer has been deleted.</li><li>• The cross-memory descriptor is not valid.</li><li>• A paging I/O error occurred while the user buffer was being accessed.</li></ul> If the user buffer is not in memory, the <code>xmemout</code> service also returns an <code>XMEM_FAIL</code> error when executing on an interrupt level.

### Related reference:

“`xmattach` Kernel Service” on page 599

“`xmemin` Kernel Service” on page 608

### Related information:

Cross Memory Kernel Services

## xmempsize Kernel Service

### Purpose

Reports the page size being used for a specified address range on the 64-bit kernel.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/xmem.h>
```

```
long long xmempsize (dp, uaddr, count)
```

```
struct xmem * dp;
void * uaddr;
size_t count;
```

### Description

The `xmempsize` kernel service returns the size, in bytes, of the virtual memory pages contained in the memory range starting at `uaddr` and continuing for `count` number of bytes. If the memory range consists of virtual memory pages of different sizes, the size of the smallest pages contained in the range is returned.

The cross-memory descriptor, `dp`, must have been previously initialized to describe the buffer containing the specified range of memory. The `xmattach()` kernel service prepares a buffer and cross-memory descriptor for use with the `xmempsize()` kernel service.

### Parameters

Item	Description
<i>dp</i>	Specifies the cross-memory descriptor.
<i>uaddr</i>	Specifies the starting address of the memory range.
<i>count</i>	Specifies the number of bytes.

## Execution Environment

The **xmempsize** kernel service can be called from either the process or interrupt environment.

The **xmempsize** kernel service is only supported on the 64-bit kernel.

## Return Values

On successful completion, the **xmempsize()** kernel service returns a page size in bytes.

Otherwise, the **xmempsize()** kernel service returns **XMEM\_FAIL**.

### Related reference:

“xmattach Kernel Service” on page 599

### Related information:

Cross Memory Kernel Services

## xmfree Kernel Service

### Purpose

Frees allocated memory.

### Syntax

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/malloc.h>
```

```
int xmfree ( ptr, heap)
caddr_t ptr;
caddr_t heap;
```

### Parameters

Item	Description
<i>ptr</i>	Specifies the address of the area in memory to free.
<i>heap</i>	Specifies the address of the heap from which the memory was allocated.

### Description

The **xmfree** kernel service frees the area of memory pointed to by the *ptr* parameter in the heap specified by the *heap* parameter. This area of memory must be allocated with the **xmalloc** kernel service. In addition, the *ptr* pointer must be the pointer returned from the corresponding **xmalloc** call.

For example, the **xmfree** (*ptr*, **kernel\_heap**) kernel service frees the area in the kernel heap allocated by *ptr=xmalloc* (*size*, *align*, **kernel\_heap**).

A kernel extension must explicitly free any memory it allocates. If it does not, eventually subsequent allocations are unsuccessful. Pinned memory must also be unpinned before it is freed if allocated from

the **kernel\_heap**. The kernel does not keep track of which kernel extension owns various allocated areas in the heap. Therefore, the kernel never automatically frees these allocated areas on process termination or device close.

An additional interface to the **xmfree** kernel service is provided. The **free** (*ptr*) is equivalent to **xmfree** (*ptr*, **kernel\_heap**).

## Execution Environment

The **xmfree** kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Indicates successful completion.
-1	Indicates one of the following errors: <ul style="list-style-type: none"><li>• The area to be freed was not allocated with the <b>xmalloc</b> kernel service.</li><li>• The heap was not initialized for memory allocation.</li></ul>

### Related reference:

“xmalloc Kernel Service” on page 598

### Related information:

Memory Kernel Services

## xmgethkeyset Kernel Service Purpose

Retrieves the hardware keyset associated with a cross-memory descriptor.

## Syntax

```
#include <sys/types.h>
#include <sys/kernno.h>
#include <sys/xmem.h>
#include <sys/skeys.h>
```

```
kernno_t xmgethkeyset (dp, keyset, flags)
struct xmem * dp;
hkeyset_t * hkeyset;
long flags;
```

## Parameters

Item	Description
<i>dp</i>	Specifies a valid cross-memory descriptor.
<i>hkeyset</i>	Pointer to returned hardware keyset associated with the cross-memory descriptor.
<i>flags</i>	Must be set to zero.

## Description

The **xmgethkeyset()** kernel service can be used to obtain the keyset associated with a cross-memory descriptor.

Kernel-key protection can be enforced on memory regions described by a cross-memory descriptor. The enforcement is done during normal access checking performed by cross-memory services, such as **xmemdma()** service.

## Execution Environment

The `xmgethkeyset` kernel service can be called from the process or interrupt environment.

## Return Values

Item	Description
0	Successful.
<code>EINVAL_XMGETHKEYSET</code>	Invalid parameter.

### Related reference:

“`xmsethkeyset` Kernel Service”

“`xmattach` Kernel Service” on page 599

## `xmsethkeyset` Kernel Service

### Purpose

Alters hardware keyset associated with a cross-memory descriptor.

### Syntax

```
#include <sys/types.h>
#include <sys/kernno.h>
#include <sys/xmem.h>
#include <sys/skeys.h>
```

```
kernno_t xmsethkeyset (dp, hkeyset, flags)
struct xmem * dp;
hkeyset_t hkeyset;
long flags;
```

### Parameters

Item	Description
<i>dp</i>	Specifies a valid cross-memory descriptor.
<i>hkeyset</i>	Hardware keyset to assign to the cross-memory descriptor.
<i>flags</i>	Must be set to zero.

### Description

The `xmsethkeyset()` kernel service can be used to modify the keyset associated with a cross-memory descriptor.

Kernel-key protection can be enforced on memory regions described by a cross-memory descriptor. The enforcement is done during normal access checking performed by cross-memory services, such as the `xmemdma()` service.

## Execution Environment

The `xmsethkeyset` kernel service can be called from the process environment only.

## Return Values

Item	Description
0	Successful.
EINVAL_XMSETHKEYSET	Invalid parameter or execution environment.

**Related reference:**

“xmgethkeyset Kernel Service” on page 612

“xmattach Kernel Service” on page 599

## Device Driver Operations

This topic provides a description of standard device driver entry points parameters.

### Standard Parameters to Device Driver Entry Points

#### Purpose

Provides a description of standard device driver entry points parameters.

#### Description

There are three parameters passed to device driver entry points that always have the same meanings: the *devno* parameter, the *chan* parameter, and the *ext* parameter.

#### The devno Parameter

This value, defined to be of type **dev\_t**, specifies the device or subdevice to which the operation is directed. For convenience and portability, the `/usr/include/sys/sysmacros.h` file defines the following macros for manipulating device numbers:

Macro	Description
<b>major</b> ( <i>devno</i> )	Returns the major device number.
<b>minor</b> ( <i>devno</i> )	Returns the minor device number.
<b>makedev</b> ( <i>maj</i> , <i>min</i> ).	Constructs a composite device number in the format of <i>devno</i> from the major and minor device numbers given.

#### The chan Parameter

This value, defined to be of type **chan\_t**, is the channel ID for a multiplexed device driver. If the device driver is not multiplexed, *chan* has the value of 0. If the driver is multiplexed, then the *chan* parameter is the **chan\_t** value returned from the device driver's **ddmpx** routine.

#### The ext Parameter

The *ext* parameter, or extension parameter, is defined to be of type **int**. It is meaningful only with calls to such extended subroutines as the **openx**, **readx**, **writex**, and **ioctlx** subroutines. These subroutines allow applications to pass an extra, device-specific parameter to the device driver. This parameter is then passed to the **ddopen**, **ddread**, **ddwrite**, and **ddioctl** device driver entry points as the *ext* parameter. If the application uses one of the non-extended subroutines (for example, the **read** instead of the **readx** subroutine), then the *ext* parameter has a value of 0.

**Note:** Using the *ext* parameter is highly discouraged because doing so makes an application program less portable to other operating systems.

**Related reference:**

“ddioctl Device Driver Entry Point” on page 626

**Related information:**

read subroutine

Device Driver Kernel Extension Overview

## buf Structure

### Purpose

Describes buffering data transfers between a program and the peripheral device

### Introduction to Kernel Buffers

For block devices, kernel buffers are used to buffer data transfers between a program and the peripheral device. These buffers are allocated in blocks of 4096 bytes. At any given time, each memory block is a member of one of two linked lists that the device driver and the kernel maintain:

#### List

**Available buffer queue (avlist)**

**Busy buffer queue (blist)**

#### Description

A list of all buffers available for use. These buffers do not contain data waiting to be transferred to or from a device.

A list of all buffers that contain data waiting to be transferred to or from a device.

Each buffer has an associated buffer header called the **buf** structure pointing to it. Each buffer header has several parts:

- Information about the block
- Flags to show status information
- Busy list forward and backward pointers
- Available list forward and backward pointers

The device driver maintains the `av_forw` and `av_back` pointers (for the available blocks), while the kernel maintains the `b_forw` and `b_back` pointers (for the busy blocks).

## buf Structure Variables for Block I/O

The **buf** structure, which is defined in the `/usr/include/sys/buf.h` file, includes the following fields:

Item	Description
b_flags	Flag bits. The value of this field is constructed by logically ORing 0 or more of the following values: <ul style="list-style-type: none"><li><b>B_WRITE</b> This operation is a write operation.</li><li><b>B_READ</b> This operation is a read data operation, rather than write.</li><li><b>B_DONE</b> I/O on the buffer has been done, so the buffer information is more current than other versions.</li><li><b>B_ERROR</b> A transfer error has occurred and the transaction has aborted.</li><li><b>B_BUSY</b> The block is not on the free list.</li><li><b>B_INFLIGHT</b> This I/O request has been sent to the physical device driver for processing.</li><li><b>B_AGE</b> The data is not likely to be reused soon, so prefer this buffer for reuse. This flag suggests that the buffer goes at the head of the free list rather than at the end.</li><li><b>B_ASYNC</b> Asynchronous I/O is being performed on this block. When I/O is done, release the block.</li><li><b>B_DELWRI</b> The contents of this buffer still need to be written out before the buffer can be reused, even though this block may be on the free list. This is used by the <b>write</b> subroutine when the system expects another write to the same block to occur soon.</li><li><b>B_NOHIDE</b> Indicates that the data page should not be hidden during direct memory access (DMA) transfer.</li><li><b>B_SETMOD</b> Allows an enhanced I/O file system to cause a page to be considered modified.</li><li><b>B_STALE</b> The data conflicts with the data on disk because of an I/O error.</li><li><b>B_XREADONLY</b> Indicates a read-only page in the external pager buffer list.</li><li><b>B_MORE_DONE</b> When set, indicates to the receiver of this <b>buf</b> structure that more structures are queued in the <b>IODONE</b> level. This permits device drivers to handle all completed requests before processing any new requests.</li><li><b>B_SPLIT</b> When set, indicates that the transfer can begin anywhere within the data buffer.</li></ul>
b_forw	The forward busy block pointer.
b_back	The backward busy block pointer.
av_forw	The forward pointer for a driver request queue.
av_back	The backward pointer for a driver request queue.
b_iodone	Anyone calling the strategy routine must set this field to point to their I/O done routine. This routine is called on the <b>INTIODONE</b> interrupt level when I/O is complete.
b_dev	The major and minor device number.
b_bcount	The byte count for the data transfer.
b_un.b_addr	The memory address of the data buffer.
b_blkno	The block number on the device.
b_resid	Amount of data not transferred after error.
b_event	Anchor for event list.
b_xmemd	Cross-memory descriptor.

### Related reference:



“ddstrategy Device Driver Entry Point” on page 635

“bufx Structure”

**Related information:**

write subroutine

Device Driver Kernel Extension Overview

## **bufx Structure**

### **Purpose**

Extends the **buf** structure to accommodate new fields as needed for performance and RAS reasons.

### **Description**

The **bufx** structure is available for use by the 64-bit kernel and 64-bit kernel extensions. The 32-bit kernel and 32-bit kernel extensions only have the option of using the **buf** structure.

## bufx Structure Variables for Block I/O

The **bufx** structure, which is defined in the `/usr/include/sys/buf.h` file, includes the following fields:

Item	Description
<code>b_flags</code>	Flag bits. The value of this field is constructed by the logical OR operation with 0 or more of the following values: <ul style="list-style-type: none"><li><b>B_WRITE</b> This operation is a write operation.</li><li><b>B_READ</b> This operation is a read data operation.</li><li><b>B_DONE</b> I/O on the buffer is done, so the buffer information is more current than other versions.</li><li><b>B_ERROR</b> A transfer error occurred and the transaction aborted.</li><li><b>B_BUSY</b> The block is not on the free list.</li><li><b>B_INFLIGHT</b> This I/O request was sent to the physical device driver for processing.</li><li><b>B_AGE</b> The data is not likely to be reused soon, so prefer this buffer for reuse. This flag suggests that the buffer goes at the head of the free list rather than at the end.</li><li><b>B_ASYNC</b> Asynchronous I/O is being performed on this block. When I/O is done, release the block.</li><li><b>B_DELWRI</b> The contents of this buffer still need to be written out before the buffer can be reused, even though this block may be on the free list. This is used by the <b>write</b> subroutine when the system expects another write to the same block to occur soon.</li><li><b>B_NOHIDE</b> Indicates that the data page should not be hidden during direct memory access (DMA) transfer.</li><li><b>B_STALE</b> The data conflicts with the data on disk because of an I/O error.</li><li><b>B_MORE_DONE</b> When set, indicates to the receiver of this <b>bufx</b> structure that more structures are queued in the <b>IODONE</b> level. This permits device drivers to handle all completed requests before processing any new requests.</li><li><b>B_SPLIT</b> When set, indicates that the transfer can begin anywhere within the data buffer.</li><li><b>B_BUFX</b> A buffer is identified as an extended <b>buf</b> structure if all of the following conditions are met:<ul style="list-style-type: none"><li><b>B_BUFX</b> bit is set in the <code>b_flags</code> field.</li><li>The pointer obtained by recombining the <code>bx_refptrtop</code> field and the <code>bx_refptrbot</code> field points to the beginning of the structure.</li><li>The <code>bx_eyecatcher</code> field, which identifies whether the <b>buf</b> structure is extended or not, is equal to the ASCII string "bufx".</li></ul></li><li><b>B_BUFX_INITIAL</b> When set, indicates that the <b>buf</b> is extended.</li></ul>
<code>b_forw</code>	The forward busy block pointer.
<code>b_back</code>	The backward busy block pointer.
<code>av_forw</code>	The forward pointer for a driver request queue.
<code>av_back</code>	The backward pointer for a driver request queue.
<code>b_iodone</code>	Anyone calling the strategy routine must set this field to point to their I/O done routine. This routine is called on the <b>INTIODONE</b> interrupt level when I/O is complete.
<code>b_dev</code>	The major and minor device number.
<code>b_bcount</code>	The byte count for the data transfer.
<code>b_un.b_addr</code>	The memory address of the data buffer.
<code>b_blkno</code>	The block number on the device.

Item	Description
b_resid	The amount of data not transferred after error.
b_event	The anchor for event list.
b_xmemd	The cross-memory descriptor.
bx_refptrtop	The top half of the reference pointer.
bx_refptrbot	The bottom half of the reference pointer.
bx_version	The version of the <b>bufx</b> structure.
bx_eyecatcher	The field contains the string "bufx", allowing for easy identification of the <b>bufx</b> structure in KDB when dumping data and for structure verification in addition to using the <b>BUFX_VALIDATE</b> macro.
bx_flags	<b>Bufx</b> flags with a 64-bit field that can be used for <b>bufx</b> -specific flags that are yet to be defined.
bx_io_priority	If the underlying storage devices do not support I/O priority, this value is ignored. The <code>bx_io_priority</code> must be either the value of <code>IOPRIORITY_UNSET</code> (0) or a value from 1 to 15. Lower I/O priority values are considered to be more important than higher values. For example, a value of 1 is considered the highest priority and a value of 15 is considered the lowest priority. The value of <code>IOPRIORITY_UNSET</code> is defined in the <code>sys/extendio.h</code> file.
bx_io_cache_hint	If the underlying storage devices do not support I/O cache hints, this value is ignored. The <code>bx_io_cache_hint</code> must be either the value of <code>CH_AGE_OUT_FAST</code> or the value of <code>CH_PAGE_WRITE</code> (defined in the <code>sys/extendio.h</code> file). These values are mutually exclusive. If <code>CH_AGE_OUT_FAST</code> is set, the I/O buffer can be aged out quickly from the storage device buffer cache. This is useful in the situations where the application is already caching the I/O buffer and redundant caching within the storage layer can be avoided. If <code>CH_PAGE_WRITE</code> is set, the I/O buffer is written only to the storage device cache and not to the disk.

#### Related reference:

“buf Structure” on page 615

## Character Lists Structure

Character device drivers, and other character-oriented support that can perform character-at-a-time I/O, can be implemented by using a common set of services and data buffers to handle characters in the form of *character lists*. A *character list* is a list or queue of characters. Some routines put characters in a list, and others remove the characters from the list.

Character lists, known as **clists**, contain a **clist** header and a chain of one or more data buffers known as character blocks. Putting characters on a queue allocates space (character blocks) from the common pool and links the character block into the data structure defining the character queue. Obtaining characters from a queue returns the corresponding space back to the pool.

A character list can be used to communicate between a character device driver top and bottom half. The **clist** header and the character blocks that are used by these routines must be pinned in memory, since they are accessed in the interrupt environment.

Users of the character list services must register (typically in the device driver **ddopen** routine) the number of character blocks to be used at any one time. This allows the kernel to manage the number of pinned character blocks in the character block pool. Similarly, when usage terminates (for example, when the device driver is closed), the using routine should remove its registration of character blocks. The **pinconf** kernel service provides registration for character block usage.

The kernel provides four services for obtaining characters or character blocks from a character list: the **getc**, **getcb**, **getcnp**, and **getcxc** kernel services. There are also four services that add characters or character blocks to character lists: the **putc**, **putcb**, **putcnp**, and **putcxc** kernel services. The **getc** kernel services allocates a free character block while the **putc** kernel service returns a character block to the free list. Additionally, the **putcfl** kernel service returns a list of character buffers to the free list. The **waitcfree** kernel service determines if any character blocks are on the free list, and waits for one if none are available.

## Using a Character List

For each character list you use, you must allocate a **clist** header structure. This **clist** structure is defined in the `/usr/include/sys/cblock.h` file.

You do not need to be concerned with maintaining the fields in the **clist** header, as the character list services do this for you. However, you should initialize the `c_cc` count field to 0, and both character block pointers (`c_cf` and `c_cl`) to null before using the **clist** header for the first time. The **clist** structure defines these fields.

Each buffer in the character list is a **cblock** structure, which is also defined in the `/usr/include/sys/cblock.h` file.

A character block data area does not need to be completely filled with characters. The `c_first` and `c_last` fields are zero-based offsets within the `c_data` array, which actually contains the data.

Only a limited amount of memory is available for character buffers. All character drivers share this pool of buffers. Therefore, you must limit the number of characters in your character list to a few hundred. When the device is closed, the device driver should make certain all of its character lists are flushed so the buffers are returned to the list of free buffers.

### Related reference:

“getc Kernel Service” on page 182

“putc Kernel Service” on page 423

### Related information:

Device Driver Kernel Extension Overview

## ddclose Device Driver Entry Point

### Purpose

Closes a previously open device instance.

### Syntax

```
#include <sys/device.h>
#include <sys/types.h>
int ddclose ( devno, chan)
dev_t devno;
chan_t chan;
```

### Parameters

Item	Description
<i>devno</i>	Specifies the major and minor device numbers of the device instance to close.
<i>chan</i>	Specifies the channel number.

### Description

The **ddclose** entry point is called when a previously opened device instance is closed by the **close** subroutine or **fp\_close** kernel service. The kernel calls the routine under different circumstances for non-multiplexed and multiplexed device drivers.

For non-multiplexed device drivers, the kernel calls the **ddclose** routine when the last process having the device instance open closes it. This causes the g-node reference count to be decremented to 0 and the g-node to be deallocated.

For multiplexed device drivers, the **ddclose** routine is called for each close associated with an explicit open. In other words, the device driver's **ddclose** routine is invoked once for each time its **ddopen** routine was invoked for the channel.

In some instances, data buffers should be written to the device before returning from the **ddclose** routine. These are buffers containing data to be written to the device that have been queued by the device driver but not yet written.

Non-multiplexed device drivers should reset the associated device to an idle state and change the device driver device state to closed. This can involve calling the **fp\_close** kernel service to issue a close to an associated open device handler for the device. Returning the device to an idle state prevents the device from generating any more interrupt or direct memory access (DMA) requests. DMA channels and interrupt levels allocated for this device should be freed, until the device is re-opened, to release critical system resources that this device uses.

Multiplexed device drivers should provide the same device quiescing, but not in the **ddclose** routine. Returning the device to the idle state and freeing its resources should be delayed until the **ddmpx** routine is called to deallocate the last channel allocated on the device.

In all cases, the device instance is considered closed once the **ddclose** routine has returned to the caller, even if a nonzero return code is returned.

## Execution Environment

The **ddclose** routine is executed only in the process environment. It should provide the required serialization of its data structures by using the locking kernel services in conjunction with a private lock word defined in the driver.

## Return Values

The **ddclose** entry point can indicate an error condition to the user-mode application program by returning a nonzero return code. This causes the subroutine call to return a value of -1. It also makes the return code available to the user-mode application in the **errno** global variable. The return code used should be one of the values defined in the **/usr/include/sys/errno.h** file.

The device is always considered closed even if a nonzero return code is returned.

When applicable, the return values defined in the POSIX 1003.1 standard for the **close** subroutine should be used.

### Related reference:

“ddopen Device Driver Entry Point” on page 629

“fp\_close Kernel Service” on page 145

### Related information:

Programming in the Kernel Environment Overview

## ddconfig Device Driver Entry Point Purpose

Performs configuration functions for a device driver.

## Syntax

```
#include <sys/device.h>
#include <sys/types.h>

int ddconfig ( devno, cmd, uiop)
dev_t devno;
int cmd;
struct uiop *uiop;
```

## Parameters

Item	Description
<i>devno</i>	Specifies the major and minor device numbers.
<i>cmd</i>	Specifies the function to be performed by the <b>ddconfig</b> routine.
<i>uiop</i>	Points to a <b>uiop</b> structure describing the relevant data area for configuration information.

## Description

The **ddconfig** entry point is used to configure a device driver. It can be called to do the following tasks:

- Initialize the device driver.
- Terminate the device driver.
- Request configuration data for the supported device.
- Perform other device-specific configuration functions.

The **ddconfig** routine is called by the device's Configure, Unconfigure, or Change method. Typically, it is called once for each device number (major and minor) to be supported. This is, however, device-dependent. The specific device method and **ddconfig** routine determines the number of times it is called.

The **ddconfig** routine can also provide additional device-specific functions relating to configuration, such as returning device vital product data (VPD). The **ddconfig** routine is usually invoked through the **sysconfig** subroutine by the device-specific Configure method.

Device drivers and their methods typically support these values for the *cmd* parameter:

Value	Description
<b>CFG_INIT</b>	Initializes the device driver and internal data areas. This typically involves the minor number specified by the <i>devno</i> parameter, for validity. The device driver's <b>ddconfig</b> routine also installs the device driver's entry points in the device switch table, if this was the first time called (for the specified major number). This can be accomplished by using the <b>devswadd</b> kernel service along with a <b>devsw</b> structure to add the device driver's entry points to the device switch table for the major device number supplied in the <i>devno</i> parameter.

The **CFG\_INIT** command parameter should also copy the device-dependent information (found in the device-dependent structure provided by the caller) into a static or dynamically allocated save area for the specified device. This information should be used when the **ddopen** routine is later called.

The device-dependent structure's address and length are described in the **uiop** structure pointed to by the *uiop* parameter. The **uiomove** kernel service can be used to copy the device-dependent structure into the device driver's data area.

When the **ddopen** routine is called, the device driver passes device-dependent information to the routines or other device drivers providing the device handler role in order to initialize the device. The delay in initializing the device until the **ddopen** call is received is useful in order to delay the use of valuable system resources (such as DMA channels and interrupt levels) until the device is actually needed.

Value	Description
<b>CFG_TERM</b>	<p>Terminates the device driver associated with the specified device number, as represented by the <i>devno</i> parameter. The <b>ddconfig</b> routine determines if any opens are outstanding on the specified <i>devno</i> parameter. If none are, the <b>CFG_TERM</b> command processing marks the device as terminated, disallowing any subsequent opens to the device. All dynamically allocated data areas associated with the specified device number should be freed.</p> <p>If this termination removes the last minor number supported by the device driver from use, the <b>devswdel</b> kernel service should be called to remove the device driver's entry points from the device switch table for the specified <i>devno</i> parameter.</p> <p>If opens are outstanding on the specified device, the terminate operation is rejected with an appropriate error code returned. The Unconfigure method can subsequently unload the device driver if all uses of it have been terminated.</p> <p>To determine if all the uses of the device driver have been terminated, a device method can make a <b>sysconfig</b> subroutine call. By using the <b>sysconfig SYS_QDVSW</b> operation, the device method can learn whether or not the device driver has removed itself from the device switch table.</p>
<b>CFG_QVPD</b>	<p>Queries device-specific vital product data (VPD).</p> <p>For this function, the calling routine sets up a <b>uio</b> structure pointed at by the <i>uiop</i> parameter to the <b>ddconfig</b> routine. This <b>uio</b> structure defines an area in the caller's storage in which the <b>ddconfig</b> routine is to write the VPD. The <b>uio</b> kernel service can be used to provide the data copy operation.</p>

The data area pointed at by the *uiop* parameter has two different purposes, depending on the *cmd* function. If the **CFG\_INIT** command has been requested, the **uio** structure describes the location and length of the device-dependent data structure (DDS) from which to read the information. If the **CFG\_QVPD** command has been requested, the **uio** structure describes the area in which to write vital product data information. The content and format of this information is established by the specific device methods in conjunction with the device driver.

The **uio** kernel service can be used to facilitate copying information into or out of this data area. The format of the **uio** structure is defined in the `/usr/include/sys/uio.h` file and described further in the **uio** structure.

## Execution Environment

The **ddconfig** routine and its operations are called in the process environment only.

## Return Values

The **ddconfig** routine sets the return code to 0 if no errors are detected for the operation specified. If an error is to be returned to the caller, a nonzero return code should be provided. The return code used should be one of the values defined in the `/usr/include/sys/errno.h` file.

If this routine was invoked by a **sysconfig** subroutine call, the return code is passed to its caller (typically a device method). It is passed by presenting the error code in the **errno** global variable and providing a -1 return code to the subroutine.

### Related reference:

“devswadd Kernel Service” on page 71

“uio

### Related information:

sysconfig subroutine

Device Driver Kernel Extension Overview

# dddump Device Driver Entry Point

## Purpose

Writes system dump data to a device.

## Syntax

```
#include <sys/device.h>
```

```
int dddump (devno, uiop, cmd, arg, chan, ext)
dev_t devno;
struct uiop * uiop;
int cmd, arg;
chan_t chan;
int ext;
```

## Parameters

Item	Description
<i>devno</i>	Specifies the major and minor device numbers.
<i>uiop</i>	Points to the <b>uiop</b> structure describing the data area or areas to be dumped.
<i>cmd</i>	The parameter from the kernel dump function that specifies the operation to be performed.
<i>arg</i>	The parameter from the caller that specifies the address of a parameter block associated with the kernel dump command.
<i>chan</i>	Specifies the channel number.
<i>ext</i>	Specifies the extension parameter.

## Description

The kernel dump routine calls the **dddump** entry point to set up and send dump requests to the device. The **dddump** routine is optional for a device driver. It is required only when the device driver supports a device as a target for a possible kernel dump.

If this is the case, it is important that the system state change as little as possible when performing the dump. As a result, the **dddump** routine should use the minimal amount of services in writing the dump data to the device.

The *cmd* parameter can specify any of the following dump commands:

Dump Command	Description
DUMPINIT	Initialization a device in preparation for supporting a system dump. The specified device instance must have previously been opened. The <i>arg</i> parameter points to a <b>dumpio_stat</b> structure, defined in <b>/usr/include/sys/dump.h</b> . This is used for returning device-specific status in case of an error.

The **dddump** routine should pin all code and data that the device driver uses to support dump writing. This is required to prevent a page fault when actually performing a write of the dump data. (Pinned code should include the **dddump** routine.) The **pin** or **pincode** kernel service can be used for this purpose.



Dump Command	Description
<b>DUMPQUERY</b>	<p>Determines the maximum and minimum number of bytes that can be transferred to the device in one <b>DUMPWRITE</b> command. For network dumps, the address of the write routine used in transferring dump data to the network dump device is also sent. The <i>uiop</i> parameter is not used and is null for this command. The <i>arg</i> parameter is a pointer to a <b>dmp_query</b> structure, as defined in the <code>/usr/include/sys/dump.h</code> file.</p> <p>The <b>dmp_query</b> structure contains the following fields:</p> <p><b>min_tsize</b> Minimum transfer size (in bytes).</p> <p><b>max_tsize</b> Maximum transfer size (in bytes).</p> <p><b>dumpwrite</b> Address of the write routine.</p> <p>The <b>DUMPQUERY</b> command returns the data transfer size information in the <b>dmp_query</b> structure pointed to by the <i>arg</i> parameter. The kernel dump function then uses a buffer between the minimum and maximum transfer sizes (inclusively) when writing dump data.</p> <p>If the buffer is not the size found in the <i>max_tsize</i> field, then its size must be a multiple of the value in the <i>min_tsize</i> field. The <i>min_tsize</i> field and the <i>max_tsize</i> field can specify the same value.</p>
<b>DUMPSTART</b>	<p>Suspends current device activity and provide whatever setup of the device is needed before receiving a <b>DUMPWRITE</b> command. The <i>arg</i> parameter points to a <b>dumpio_stat</b> structure, defined in <code>/usr/include/sys/dump.h</code>. This is used for returning device-specific status in case of an error.</p>
<b>DUMPWRITE</b>	<p>Writes dump data to the target device. The <b>uio</b> structure pointed to by the <i>uiop</i> parameter specifies the data area or areas to be written to the device and the starting device offset. The <i>arg</i> parameter points to a <b>dumpio_stat</b> structure, defined in <code>/usr/include/sys/dump.h</code>. This is used for returning device-specific status in case of an error. Code for the <b>DUMPWRITE</b> command should minimize its reliance on system services, process dispatching, and such interrupt services as the <b>INTIODONE</b> interrupt priority or device hardware interrupts.</p> <p><b>Note:</b> The <b>DUMPWRITE</b> command must never cause a page fault. This is ensured on the part of the caller, since the data areas to be dumped have been determined to be in memory. The device driver must ensure that all of its code, data and stack accesses are to pinned memory during its <b>DUMPINIT</b> command processing.</p>
<b>DUMPEND</b>	<p>Indicates that the kernel dump has been completed. Any cleanup of the device state should be done at this time.</p>
<b>DUMPTERM</b>	<p>Indicates that the specified device is no longer a selected dump target device. If no other devices supported by this <b>dddump</b> routine have a <b>DUMPINIT</b> command outstanding, the <b>DUMPTERM</b> code should unpin any resources pinned when it received the <b>DUMPINIT</b> command. (The <b>unpin</b> kernel service is available for unpinning memory.) The <b>DUMPTERM</b> command is received before the device is closed.</p>
<b>DUMPREAD</b>	<p>Receives the acknowledgment packet for previous <b>DUMPWRITE</b> operations to a communications device driver. If the device driver receives the acknowledgment within the specified time, it returns a 0 and the response data is returned to the kernel dump function in the <i>uiop</i> parameter. If the device driver does not receive the acknowledgment within the specified time, it returns a value of <b>ETIMEDOUT</b>.</p> <p>The <i>arg</i> parameter contains a timeout value in milliseconds.</p>

## Execution Environment

The **DUMPINIT dddump** operation is called in the process environment only. The **DUMPQUERY**, **DUMPSTART**, **DUMPWRITE**, **DUMPEND**, and **DUMPTERM dddump** operations can be called in both the process environment and interrupt environment.

## Return Values

The **dddump** entry point indicates an error condition to the caller by returning a nonzero return code.

### Related reference:

“devdump Kernel Service” on page 69

“dmp\_add Kernel Service” on page 92

### Related information:

Device Driver Kernel Extension Overview

# ddioctl Device Driver Entry Point

## Purpose

Performs the special I/O operations requested in an **ioctl** or **ioctlx** subroutine call.

## Syntax

```
#include <sys/device.h>
```

```
int ddiioctl (devno, cmd, arg, devflag, chan, ext)
dev_t devno;
int cmd;
void *arg;
ulong devflag;
chan_t chan;
int ext;
```

## Description

When a program issues an **ioctl** or **ioctlx** subroutine call, the kernel calls the **ddioctl** routine of the specified device driver. The **ddioctl** routine is responsible for performing whatever functions are requested. In addition, it must return whatever control information has been specified by the original caller of the **ioctl** subroutine. The *cmd* parameter contains the name of the operation to be performed.

Most **ioctl** operations depend on the specific device involved. However, all **ioctl** routines must respond to the following command:

Item	Description
IOCINFO	Returns a <b>devinfo</b> structure (defined in the <code>/usr/include/sys/devinfo.h</code> file) that describes the device. (Refer to the description of the special file for a particular device in the Application Programming Interface.) Only the first two fields of the data structure need to be returned if the remaining fields of the structure do not apply to the device.

The *devflag* parameter indicates one of several types of information. It can give conditions in which the device was opened. (These conditions can subsequently be changed by the **fcntl** subroutine call.) Alternatively, it can tell which of two ways the entry point was invoked:

- By the file system on behalf of a using application
- Directly by a kernel routine using the **fp\_ioctl** kernel service

Thus flags in the *devflag* parameter have the following definitions, as defined in the `/usr/include/sys/device.h` file:

Item	Description
DKERNEL	Entry point called by kernel routine using the <b>fp_ioctl</b> service.
DREAD	Open for reading.
DWRITE	Open for writing.
DAPPEND	Open for appending.
DNDELAY	Device open in nonblocking mode.

## Parameters

Item	Description
<i>devno</i>	Specifies the major and minor device numbers.
<i>cmd</i>	The parameter from the <b>ioctl</b> subroutine call that specifies the operation to be performed.
<i>arg</i>	The parameter from the <b>ioctl</b> subroutine call that specifies an additional argument for the <i>cmd</i> operation.
<i>devflag</i>	Specifies the device open or file control flags.
<i>chan</i>	Specifies the channel number.
<i>ext</i>	Specifies the extension parameter.

## Execution Environment

The **ddioctl** routine is executed only in the process environment. It should provide the required serialization of its data structures by using the locking kernel services in conjunction with a private lock word defined in the driver.

## Return Values

The **ddioctl** entry point can indicate an error condition to the user-mode application program by returning a nonzero return code. This causes the **ioctl** subroutine to return a value of -1 and makes the return code available to the user-mode application in the **errno** global variable. The error code used should be one of the values defined in the `/usr/include/sys/errno.h` file.

When applicable, the return values defined in the POSIX 1003.1 standard for the **ioctl** subroutine should be used.

### Related reference:

“Standard Parameters to Device Driver Entry Points” on page 614

“fp\_ioctl Kernel Service” on page 151

### Related information:

fcntl subroutine

Device Driver Kernel Extension Overview

## ddmpx Device Driver Entry Point

### Purpose

Allocates or deallocates a channel for a multiplexed device driver.

### Syntax

```
#include <sys/device.h>
#include <sys/types.h>
```

```
int ddmpx ( devno, chanp, channame)
dev_t devno;
chan_t *chanp;
char *channame;
```

### Parameters

Item	Description
<i>devno</i>	Specifies the major and minor device numbers.
<i>chanp</i>	Specifies the channel ID, passed by reference.
<i>channame</i>	Points to the path name extension for the channel to be allocated.

## Description

Only multiplexed character class device drivers can provide the **ddmpx** routine, and *every* multiplexed driver must do so. The **ddmpx** routine cannot be provided by block device drivers even when providing *raw* read/write access.

A multiplexed device driver is a character class device driver that supports the assignment of channels to provide finer access control to a device or virtual subdevice. This type of device driver has the capability to decode special channel-related information appended to the end of the path name of the device's special file. This path name extension is used to identify a logical or virtual subdevice or channel.

When an **open** or **creat** subroutine call is issued to a device instance supported by a multiplexed device driver, the kernel calls the device driver's **ddmpx** routine to allocate a channel.

The kernel calls the **ddmpx** routine when a channel is to be allocated or deallocated. Upon allocation, the kernel dynamically creates g-nodes (in-core i-nodes) for channels on a multiplexed device to allow the protection attributes to differ for various channels.

To allocate a channel, the **ddmpx** routine is called with a *channame* pointer to the path name extension. The path name extension starts after the first / (slash) character that follows the special file name in the path name. The **ddmpx** routine should perform the following actions:

- Parse this path name extension.
- Allocate the corresponding channel.
- Return the channel ID through the *chanp* parameter.

If no path name extension exists, the *channame* pointer points to a null character string. In this case, an available channel should be allocated and its channel ID returned through the *chanp* parameter.

If no error is returned from the **ddmpx** routine, the returned channel ID is used to determine if the channel was already allocated. If already allocated, the g-node for the associated channel has its reference count incremented. If the channel was not already allocated, a new g-node is created for the channel. In either case, the device driver's **ddopen** routine is called with the channel number assigned by the **ddmpx** routine. If a nonzero return code is returned by the **ddmpx** routine, the channel is assumed not to have been allocated, and the device driver's **ddopen** routine is not called.

If a close of a channel is requested so that the channel is no longer used (as determined by the channel's g-node reference count going to 0), the kernel calls the **ddmpx** routine. The **ddmpx** routine deallocates the channel after the **ddclose** routine was called to close the last use of the channel. If a nonzero return code is returned by the **ddclose** routine, the **ddmpx** routine is still called to deallocate the channel. The **ddclose** routine's return code is saved, to be returned to the caller. If the **ddclose** routine returned no error, but a nonzero return code was returned by the **ddmpx** routine, the channel is assumed to be deallocated, although the return code is returned to the caller.

To deallocate a channel, the **ddmpx** routine is called with a null *channame* pointer and the channel ID passed by reference in the *chanp* parameter. If the channel g-node reference count has gone to 0, the kernel calls the **ddmpx** routine to deallocate the channel after invoking the **ddclose** routine to close it. The **ddclose** routine should not itself deallocate the channel.

## Execution Environment

The **ddmpx** routine is called in the process environment only.

## Return Values

If the allocation or deallocation of a channel is successful, the **ddmpx** routine should return a return code of 0. If an error occurs on allocation or deallocation, this routine returns a nonzero value.

The return code should conform to the return codes described for the **open** and **close** subroutines in the POSIX 1003.1 standard, where applicable. Otherwise, the return code should be one defined in the `/usr/include/sys/errno.h` file.

### Related reference:

“ddclose Device Driver Entry Point” on page 620

“ddopen Device Driver Entry Point”

### Related information:

Device Driver Kernel Extension Overview

## ddopen Device Driver Entry Point

### Purpose

Prepares a device for reading, writing, or control functions.

### Syntax

```
#include <sys/device.h>
int ddopen (devno, devflag, chan, ext)
dev_t devno;
ulong devflag;
chan_t chan;
int ext;
```

### Parameters

Item	Description
<i>devno</i>	Indicates major and minor device numbers.
<i>devflag</i>	Specifies open file control flags.
<i>chan</i>	Specifies the channel number.
<i>ext</i>	Specifies the extension parameter.

### Description

The kernel calls the **ddopen** routine of a device driver when a program issues an **open** or **creat** subroutine call. It can also be called when a system call, kernel process, or other device driver uses the **fp\_opendev** or **fp\_open** kernel service to use the device.

The **ddopen** routine must first ensure exclusive access to the device, if necessary. Many character devices, such as printers and plotters, should be opened by only one process at a time. The **ddopen** routine can enforce this by maintaining a static flag variable, which is set to 1 if the device is open and 0 if not.

Each time the **ddopen** routine is called, it checks the value of the flag. If the value is other than 0, the **ddopen** routine returns with a return code of **EBUSY** to indicate that the device is already open. Otherwise, the **ddopen** routine sets the flag and returns normally. The **ddclose** entry point later clears the flag when the device is closed.

Since most block devices can be used by several processes at once, a block driver should not try to enforce opening by a single user.

The **ddopen** routine must initialize the device if this is the first open that has occurred. Initialization involves the following steps:

1. The **ddopen** routine should allocate the required system resources to the device (such as DMA channels, interrupt levels, and priorities). It should, if necessary, register its device interrupt handler for the interrupt level required to support the target device. (The **i\_init** and **d\_init** kernel services are available for initializing these resources.)
2. If this device driver is providing the head role for a device and another device driver is providing the handler role, the **ddopen** routine should use the **fp\_opendev** kernel service to open the device handler.

**Note:** The **fp\_opendev** kernel service requires a *devno* parameter to identify which device handler to open. This *devno* value, taken from the appropriate device dependent structure (DDS), should have been stored in a special save area when this device driver's **ddconfig** routine was called.

### Flags Defined for the devflag Parameter

The *devflag* parameter has the following flags, as defined in the `/usr/include/sys/device.h` file:

Item	Description
<b>DKERNEL</b>	Entry point called by kernel routine using the <b>fp_opendev</b> or <b>fp_open</b> kernel service.
<b>DREAD</b>	Open for reading.
<b>DWRITE</b>	Open for writing.
<b>DAPPEND</b>	Open for appending.
<b>DNDELAY</b>	Device open in nonblocking mode.

### Execution Environment

The **ddopen** routine is executed only in the process environment. It should provide the required serialization of its data structures by using the locking kernel services in conjunction with a private lock word defined in the driver.

### Return Values

The **ddopen** entry point can indicate an error condition to the user-mode application program by returning a nonzero return code. Returning a nonzero return code causes the **open** or **creat** subroutines to return a value of -1 and makes the return code available to the user-mode application in the **errno** global variable. The return code used should be one of the values defined in the `/usr/include/errno.h` file.

If a nonzero return code is returned by the **ddopen** routine, the open request is considered to have failed. No access to the device instance is available to the caller as a result. In addition, for nonmultiplexed drivers, if the failed open was the first open of the device instance, the kernel calls the driver's **ddclose** entry point to allow resources and device driver state to be cleaned up. If the driver was multiplexed, the kernel does not call the **ddclose** entry point on an open failure.

When applicable, the return values defined in the POSIX 1003.1 standard for the **open** subroutine should be used.

#### Related reference:

“ddclose Device Driver Entry Point” on page 620

#### Related information:

close subroutine

Programming in the Kernel Environment Overview

## ddread Device Driver Entry Point

### Purpose

Reads in data from a character device.

### Syntax

```
#include <sys/device.h>
#include <sys/types.h>

int ddread ( devno, uiop, chan, ext)
dev_t devno;
struct uio *uiop;
chan_t chan;
int ext;
```

### Parameters

Item	Description
<i>devno</i>	Specifies the major and minor device numbers.
<i>uiop</i>	Points to a <b>uio</b> structure describing the data area or areas in which to be written.
<i>chan</i>	Specifies the channel number.
<i>ext</i>	Specifies the extension parameter.

### Description

When a program issues a **read** or **readx** subroutine call or when the **fp\_rwuio** kernel service is used, the kernel calls the **ddread** entry point.

This entry point receives a pointer to a **uio** structure that provides variables used to specify the data transfer operation.

Character device drivers can use the **ureadc** and **uiomove** kernel services to transfer data into and out of the user buffer area during a **read** subroutine call. These services receive a pointer to the **uio** structure and update the fields in the structure by the number of bytes transferred. The only fields in the **uio** structure that cannot be modified by the data transfer are the **uio\_fmode** and **uio\_segflg** fields.

For most devices, the **ddread** routine sends the request to the device handler and then waits for it to finish. The waiting can be accomplished by calling the **e\_sleep** kernel service. This service suspends the driver and the process that called it and permits other processes to run until a specified event occurs.

When the I/O operation completes, the device usually issues an interrupt, causing the device driver's interrupt handler to be called. The interrupt handler then calls the **e\_wakeup** kernel service specifying the awaited event, thus allowing the **ddread** routine to resume.

The **uio\_resid** field initially contains the total number of bytes to read from the device. If the device driver supports it, the **uio\_offset** field indicates the byte offset on the device from which the read should start.

The **uio\_offset** field is a 64 bit integer (**offset\_t**); this allows the file system to send I/O requests to a device driver's read & write entry points which have logical offsets beyond 2 gigabytes. Device drivers must use care not to cause a loss of significance by assigning the offset to a 32 bit variable or using it in calculations that overflow a 32 bit variable.

If no error occurs, the `uio_resid` field should be 0 on return from the **ddread** routine to indicate that all requested bytes were read. If an error occurs, this field should contain the number of bytes remaining to be read when the error occurred.

If a read request starts at a valid device offset but extends past the end of the device's capabilities, no error should be returned. However, the `uio_resid` field should indicate the number of bytes not transferred. If the read starts at the end of the device's capabilities, no error should be returned. However, the `uio_resid` field should not be modified, indicating that no bytes were transferred. If the read starts past the end of the device's capabilities, an **ENXIO** return code should be returned, without modifying the `uio_resid` field.

When the **ddread** entry point is provided for raw I/O to a block device, this routine usually translates requests into block I/O requests using the **uphysio** kernel service.

## Execution Environment

The **ddread** routine is executed only in the process environment. It should provide the required serialization of its data structures by using the locking kernel services in conjunction with a private lock word defined in the driver.

## Return Values

The **ddread** entry point can indicate an error condition to the caller by returning a nonzero return code. This causes the subroutine call to return a value of -1. It also makes the return code available to the user-mode program in the **errno** global variable. The error code used should be one of the values defined in the `/usr/include/sys/errno.h` file.

When applicable, the return values defined in the POSIX 1003.1 standard for the **read** subroutine should be used.

### Related reference:

“ddwrite Device Driver Entry Point” on page 637

“Select/Poll Logic for ddwrite and ddread Routines” on page 639

### Related information:

`read`, `readx`

Programming in the Kernel Environment Overview

## ddrevoke Device Driver Entry Point Purpose

Ensures that a secure path to a terminal is provided.

## Syntax

```
#include <sys/device.h>
```

```
#include <sys/types.h>
```

```
int ddrevoke ( devno, chan, flag)
```

```
dev_t devno;
```

```
chan_t chan;
```

```
int flag;
```

## Parameters



Item	Description
<i>devno</i>	Specifies the major and minor device numbers.
<i>chan</i>	Specifies the channel number. For a multiplexed device driver, a value of -1 in this parameter means access to all channels is to be revoked.
<i>flag</i>	Currently defined to have the value of 0. (Reserved for future extensions.)

## Description

The **ddrevoke** entry point can be provided only by character class device drivers. It cannot be provided by block device drivers even when providing raw read/write access. A **ddrevoke** entry point is required only by device drivers supporting devices in the Trusted Computing Path to a terminal (for example, by the `/dev/lft` and `/dev/tty` files for the low function terminal and teletype device drivers). The **ddrevoke** routine is called by the **frevoke** and **revoke** subroutines.

The **ddrevoke** routine revokes access to a specific device or channel (if the device driver is multiplexed). When called, the **ddrevoke** routine should terminate all processes waiting in the device driver while accessing the specified device or channel. It should terminate the processes by sending a SIGKILL signal to all processes currently waiting for a specified device or channel data transfer. The current process is not to be terminated.

If the device driver is multiplexed and the channel ID in the *chan* parameter has the value -1, all channels are to be revoked.

## Execution Environment

The **ddrevoke** routine is called in the process environment only.

## Return Values

The **ddrevoke** routine should return a value of 0 for successful completion, or a value from the `/usr/include/errno.h` file on error.

## Files

Item	Description
<code>/dev/lft</code>	Specifies the path of the LFT special file.
<code>/dev/tty</code>	Specifies the path of the tty special file.

### Related information:

frevoke subroutine

revoke subroutine

TTY Subsystem Overview

## ddselect Device Driver Entry Point

### Purpose

Checks to see if one or more events has occurred on the device.

### Syntax

```
#include <sys/device.h>
#include <sys/poll.h>
```

```
int ddselect ( devno, events, reventp, chan)
dev_t devno;
```

```

ushort events;
ushort *reventp;
int chan;

```

## Parameters

Item	Description
<i>devno</i>	Specifies the major and minor device numbers.
<i>events</i>	Specifies the events to be checked.
<i>reventp</i>	Returned events pointer. This parameter, passed by reference, is used by the <b>ddselect</b> routine to indicate which of the selected events are true at the time of the call. The returned events location pointed to by the <i>reventp</i> parameter is set to 0 before entering this routine.
<i>chan</i>	Specifies the channel number.

## Description

The **ddselect** entry point is called when the **select** or **poll** subroutine is used, or when the **fp\_select** kernel service is invoked. It determines whether a specified event or events have occurred on the device.

Only character class device drivers can provide the **ddselect** routine. It cannot be provided by block device drivers even when providing raw read/write access.

### Requests for Information on Events

The *events* parameter represents possible events to check as flags (bits). There are three basic events defined for the **select** and **poll** subroutines, when applied to devices supporting select or poll operations:

Event	Description
<b>POLLIN</b>	Input is present on the device.
<b>POLLOUT</b>	The device is capable of output.
<b>POLLPRI</b>	An exceptional condition has occurred on the device.

A fourth event flag is used to indicate whether the **ddselect** routine should record this request for later notification of the event using the **selnotify** kernel service. This flag can be set in the *events* parameter if the device driver is not required to provide asynchronous notification of the requested events:

Event	Description
<b>POLLSYNC</b>	This request is a synchronous request only. The routine need not call the <b>selnotify</b> kernel service for this request even if the events later occur.

Additional event flags in the *events* parameter are left for device-specific events on the **poll** subroutine call.

### Select Processing

If one or more events specified in the *events* parameter are true, the **ddselect** routine should indicate this by setting the corresponding bits in the *reventp* parameter. Note that the *reventp* returned events parameter is passed by reference.

If none of the requested events are true, then the **ddselect** routine sets the returned events parameter to 0. It is passed by reference through the *reventp* parameter. It also checks the **POLLSYNC** flag in the *events* parameter. If this flag is true, the **ddselect** routine should just return, since the event request was a synchronous request only.

However, if the **POLLSYNC** flag is false, the **ddselect** routine must notify the kernel when one or more of the specified events later happen. For this purpose, the routine should set separate internal flags for each event requested in the *events* parameter.

When any of these events become true, the device driver routine should use the **selnotify** service to notify the kernel. The corresponding internal flags should then be reset to prevent re-notification of the event.

Sometimes the device can be in a state in which a supported event or events can never be satisfied (such as when a communication line is not operational). In this case, the **ddselect** routine should simply set the corresponding *reventp* flags to 1. This prevents the **select** or **poll** subroutine from waiting indefinitely. As a result however, the caller will not in this case be able to distinguish between satisfied events and unsatisfiable ones. Only when a later request with an **NDELAY** option fails will the error be detected.

**Note:** Other device driver routines (such as the **ddread**, **ddwrite** routines) may require logic to support select or poll operations.

## Execution Environment

The **ddselect** routine is executed only in the process environment. It should provide the required serialization of its data structures by using the locking kernel services in conjunction with a private lock word defined in the driver.

## Return Values

The **ddselect** routine should return with a return code of 0 if the select or poll operation requested is valid for the resource specified. Requested operations are not valid, however, if either of the following is true:

- The device driver does not support a requested event.
- The device is in a state in which poll and select operations are not accepted.

In these cases, the **ddselect** routine should return with a nonzero return code (typically **EINVAL**), and without setting the relevant *reventp* flags to 1. This causes the **poll** subroutine to return to the caller with the **POLLERR** flag set in the returned events parameter associated with this resource. The **select** subroutine indicates to the caller that all requested events are true for this resource.

When applicable, the return values defined in the POSIX 1003.1 standard for the **select** subroutine should be used.

### Related reference:

“fp\_select Kernel Service” on page 167

### Related information:

select subroutine

Programming in the Kernel Environment Overview

## ddstrategy Device Driver Entry Point

### Purpose

Performs block-oriented I/O by scheduling a read or write to a block device.

### Syntax

```
void ddstrategy ( bp)
struct buf *bp;
```

## Parameter

Item	Description
<i>bp</i>	Points to a <b>buf</b> structure describing all information needed to perform the data transfer.

## Description

When the kernel needs a block I/O transfer, it calls the **ddstrategy** strategy routine of the device driver for that device. The strategy routine schedules the I/O to the device. This typically requires the following actions:

- The request or requests must be added on the list of I/O requests that need to be processed by the device.
- If the request list was empty before the preceding additions, the device's start I/O routine must be called.

## Required Processing

The **ddstrategy** routine can receive a single request with multiple **buf** structures. However, it is not required to process requests in any specific order.

The strategy routine can be passed a list of operations to perform. The *av\_forw* field in the **buf** header describes this null-terminated list of **buf** headers. This list is not doubly linked: the *av\_back* field is undefined.

Block device drivers must be able to perform multiple block transfers. If the device cannot do multiple block transfers, or can only do multiple block transfers under certain conditions, then the device driver must transfer the data with more than one device operation.

## Kernel Buffers and Using the buf Structure

An area of memory is set aside within the kernel memory space for buffering data transfers between a program and the peripheral device. Each kernel buffer has a header, the **buf** structure, which contains all necessary information for performing the data transfer. The **ddstrategy** routine is responsible for updating fields in this header as part of the transfer.

The caller of the strategy routine should set the *b\_iodone* field to point to the caller's I/O done routine. When an I/O operation is complete, the device driver calls the **iodone** kernel service, which then calls the I/O done routine specified in the *b\_iodone* field. The **iodone** kernel service makes this call from the **INTIODONE** interrupt level.

The value of the *b\_flags* field is constructed by logically ORing zero or more possible *b\_flags* field flag values.

### Attention:

- Do not modify any of the following fields of the **buf** structure passed to the **ddstrategy** entry point: the *b\_forw*, *b\_back*, *b\_dev*, *b\_un*, or *b\_blkno* field. Modifying these fields can cause unpredictable and disastrous results.
- Do not modify any of the following fields of a **buf** structure acquired with the **geteblk** service: the *b\_flags*, *b\_forw*, *b\_back*, *b\_dev*, *b\_count*, or *b\_un* field. Modifying any of these fields can cause unpredictable and disastrous results.

## Execution Environment

The **ddstrategy** routine must be coded to execute in an interrupt handler execution environment (device driver bottom half). That is, the routine should neither touch user storage, nor page fault, nor sleep.

## Return Values

The **ddstrategy** routine, unlike other device driver routines, does not return a return code. Any error information is returned in the appropriate fields within the **buf** structure pointed to by the *bp* parameter.

When applicable, the return values defined in the POSIX 1003.1 standard for the **read** and **write** subroutines must be used.

### Related reference:

“getblk Kernel Service” on page 186

### Related information:

read subroutine

write subroutine

## ddwrite Device Driver Entry Point

### Purpose

Writes out data to a character device.

### Syntax

```
#include <sys/device.h>
#include <sys/types.h>
```

```
int ddrive (devno, uiop, chan, ext)
dev_t devno;
struct uio * uiop;
chan_t chan;
int ext;
```

### Parameters

Item	Description
<i>devno</i>	Specifies the major and minor device numbers.
<i>uiop</i>	Points to a <b>uio</b> structure describing the data area or areas from which to be written.
<i>chan</i>	Specifies the channel number.
<i>ext</i>	Specifies the extension parameter.

### Description

When a program issues a **write** or **writex** subroutine call or when the **fp\_rwuio** kernel service is used, the kernel calls the **ddwrite** entry point.

This entry point receives a pointer to a **uio** structure, which provides variables used to specify the data transfer operation.

Character device drivers can use the **uwritec** and **uiomove** kernel services to transfer data into and out of the user buffer area during a **write** subroutine call. These services are passed a pointer to the **uio** structure. They update the fields in the structure by the number of bytes transferred. The only fields in the **uio** structure that are not potentially modified by the data transfer are the **uio\_fmode** and **uio\_segflg** fields.

For most devices, the **ddwrite** routine queues the request to the device handler and then waits for it to finish. The waiting is typically accomplished by calling the **e\_sleep** kernel service to wait for an event. The **e\_sleep** kernel service suspends the driver and the process that called it and permits other processes to run.

When the I/O operation is completed, the device usually causes an interrupt, causing the device driver's interrupt handler to be called. The interrupt handler then calls the **e\_wakeup** kernel service specifying the awaited event, thus allowing the **ddwrite** routine to resume.

The **uio\_resid** field initially contains the total number of bytes to write to the device. If the device driver supports it, the **uio\_offset** field indicates the byte offset on the device from where the write should start.

The **uio\_offset** field is a 64 bit integer (**offset\_t**); this allows the file system to send I/O requests to a device driver's read & write entry points which have logical offsets beyond 2 gigabytes. Device drivers must use care not to cause a loss of significance by assigning the offset to a 32 bit variable or using it in calculations that overflow a 32 bit variable.

If no error occurs, the **uio\_resid** field should be 0 on return from the **ddwrite** routine to indicate that all requested bytes were written. If an error occurs, this field should contain the number of bytes remaining to be written when the error occurred.

If a write request starts at a valid device offset but extends past the end of the device's capabilities, no error should be returned. However, the **uio\_resid** field should indicate the number of bytes not transferred. If the write starts at or past the end of the device's capabilities, no data should be transferred. An error code of **ENXIO** should be returned, and the **uio\_resid** field should not be modified.

When the **ddwrite** entry point is provided for raw I/O to a block device, this routine usually uses the **uphysio** kernel service to translate requests into block I/O requests.

## Execution Environment

The **ddwrite** routine is executed only in the process environment. It should provide the required serialization of its data structures by using the locking kernel services in conjunction with a private lock word defined in the driver.

## Return Values

The **ddwrite** entry point can indicate an error condition to the caller by returning a nonzero return value. This causes the subroutine to return a value of -1. It also makes the return code available to the user-mode program in the **errno** global variable. The error code used should be one of the values defined in the **/usr/include/sys/errno.h** file.

When applicable, the return values defined in the POSIX 1003.1 standard for the **write** subroutine should be used.

## Related Information

The **ddread** device driver entry point.

The **CIO\_GET\_FASTWRT** **ddioctl**.

The **e\_sleep** kernel service, **e\_wakeup** kernel service, **fp\_rwuio** kernel service, **uio** kernel service, **uphysio** kernel service, **uwritec** kernel service.

The **uio** structure.

The **write** and **writex** subroutines.

Device Driver Kernel Extension Overview, Understanding Device Driver Roles, Understanding Interrupts, Understanding Locking in *Kernel Extensions and Device Support Programming Concepts*.

**Related reference:**

“ddread Device Driver Entry Point” on page 631

**Related information:**

CIO\_GET\_FASTWRT subroutine

Device Driver Kernel Extension Overview

## Select/Poll Logic for ddwrite and ddread Routines

### Description

The **ddread** and **ddwrite** entry points require logic to support the **select** and **poll** operations. Depending on how the device driver is written, the interrupt routine may also need to include this logic as well.

The select/poll logic is required wherever code checks on the occurrence of desired events. At each point where one of the selection criteria is found to be true, the device driver should check whether a notification is due for that selection. If so, it should call the **selnotify** kernel service to notify the kernel of the event.

The *devno*, *chan*, and *revents* parameters are passed to the **selnotify** kernel service to indicate which device and which events have become true.

**Related reference:**

“ddread Device Driver Entry Point” on page 631

**Related information:**

poll subroutine

Device Driver Kernel Extension Overview

## uio Structure

### Purpose

Describes a memory buffer to be used in a data transfer.

### Introduction

The user I/O or **uio** structure is a data structure describing a memory buffer to be used in a data transfer. The **uio** structure is most commonly used in the read and write interfaces to device drivers supporting character or raw I/O. It is also useful in other instances in which an input or output buffer can exist in different kinds of address spaces, and in which the buffer is not contiguous in virtual memory.

The **uio** structure is defined in the `/usr/include/sys/uio.h` file.

### Description

The **uio** structure describes a buffer that is not contiguous in virtual memory. It also indicates the address space in which the buffer is defined. When used in the character device read and write interface, it also contains the device open-mode flags, along with the device read/write offset.

The kernel provides services that access data using a **uio** structure. The **ureadc**, **uwritec**, **uimove**, and **uphysio** kernel services all perform data transfers into or out of a data buffer described by a **uio** structure. The **ureadc** kernel service writes a character into the buffer described by the **uio** structure. The **uwritec** kernel service reads a character from the buffer. These two services have names opposite from what you would expect, since they are named for the user action initiating the operation. A read on the part of the user thus results in a device driver writing to the buffer, while a write results in a driver reading from the buffer.

The **uiomove** kernel service copies data to or from a buffer described by a **uio** structure from or to a buffer in the system address space. The **uphysio** kernel service is used primarily by block device drivers providing raw I/O support. The **uphysio** kernel service converts the character read or write request into a block read or write request and sends it to the **ddstrategy** routine.

The buffer described by the **uio** structure can consist of multiple noncontiguous areas of virtual memory of different lengths. This is achieved by describing the data buffer with an array of elements, each of which consists of a virtual memory address and a byte length. Each element is defined as an **iovec** element. The **uio** structure also contains a field specifying the total number of bytes in the data buffer described by the structure.

Another field in the **uio** structure describes the address space of the data buffer, which can either be system space, user space, or cross-memory space. If the address space is defined as cross memory, an additional array of cross-memory descriptors is specified in the **uio** structure to match the array of **iovec** elements.

The **uio** structure also contains a byte offset (**uio\_offset**). This field is a 64 bit integer (**offset\_t**); it allows the file system to send I/O requests to a device driver's read & write entry points which have logical offsets beyond 2 gigabytes. Device drivers must use care not to cause a loss of significance by assigning the offset to a 32 bit variable or using it in calculations that overflow a 32 bit variable.

The called routine (device driver) is permitted to modify fields in the **uio** and **iovec** structures as the data transfer progresses. The final **uio\_resid** count is in fact used to determine how much data was transferred. Therefore this count must be decremented, with each operation, by the number of bytes actually copied.

The **uio** structure contains the following fields:

Field	Description
<b>uio_iov</b>	A pointer to an array of <b>iovec</b> structures describing the user buffer for the data transfer.
<b>uio_xmem</b>	A pointer to an array of <b>xmem</b> structures containing the cross-memory descriptors for the <b>iovec</b> array.
<b>uio_iovcnt</b>	The number of yet-to-be-processed <b>iovec</b> structures in the array pointed to by the <b>uio_iov</b> pointer. The count must be at least 1. If the count is greater than 1, then a <i>scatter-gather</i> of the data is to be performed into or out of the areas described by the <b>iovec</b> structures.
<b>uio_iovdcnt</b>	The number of already processed <b>iovec</b> structures in the <b>iovec</b> array.
<b>uio_offset</b>	The file offset established by a previous <b>lseek</b> , <b>llseek</b> subroutine call. Most character devices ignore this variable, but some, such as the <b>/dev/mem</b> pseudo-device, use and maintain it.
<b>uio_segflg</b>	A flag indicating the type of buffer being described by the <b>uio</b> structure. This flag typically describes whether the data area is in user or kernel space or is in cross-memory. Refer to the <b>/usr/include/sys/uio.h</b> file for a description of the possible values of this flag and their meanings.
<b>uio_fmode</b>	The value of the file mode that was specified on opening the file or modified by the <b>fcntl</b> subroutine. This flag describes the file control parameters. The <b>/usr/include/sys/fcntl.h</b> file contains specific values for this flag.
<b>uio_resid</b>	The byte count for the data transfer. It must not exceed the sum of all the <b>iov_len</b> values in the array of <b>iovec</b> structures. Initially, this field contains the total byte count, and when the operation completes, the value must be decremented by the actual number of bytes transferred.

The **iovec** structure contains the starting address and length of a contiguous data area to be used in a data transfer. The **iovec** structure is the element type in an array pointed to by the **uio\_iov** field in the **uio** structure. This array can contain any number of **iovec** structures, each of which describes a single unit of contiguous storage. Taken together, these units represent the total area into which, or from which, data is to be transferred. The **uio\_iovcnt** field gives the number of **iovec** structures in the array.

The **iovec** structure contains the following fields:



Field	Description
iov_base	A variable in the <b>iovec</b> structure containing the base address of the contiguous data area in the address space specified by the <code>uio_segflag</code> field. The length of the contiguous data area is specified by the <code>iov_len</code> field.
iov_len	A variable in the <b>iovec</b> structure containing the byte length of the data area starting at the address given in the <code>iov_base</code> variable.

**Related reference:**

“uio\_move Kernel Service” on page 517

“uphysio Kernel Service” on page 524

“vnode\_getacl Entry Point” on page 664

**Related information:**

Device Driver Kernel Extension Overview

## Virtual File System Operations

The following topic provides entry points specified by the virtual file system interface for performing operations on `vfs` structures.

The following entry points are specified by the virtual file system interface for performing operations on `vfs` structures:

Entry Point	Description
<code>vfs_aclxctl</code>	Issues ACL related control operations for a file system.
<code>vfs_ctl</code>	Issues control operations for a file system.
<code>vfs_init</code>	Initializes a virtual file system.
<code>vfs_mount</code>	Mounts a virtual file system.
<code>vfs_root</code>	Finds the root v-node of a virtual file system.
<code>vfs_statfs</code>	Obtains virtual file system statistics.
<code>vfs_sync</code>	Forces file system updates to permanent storage.
<code>vfs_umount</code>	Unmounts a virtual file system.
<code>vfs_vget</code>	Gets the v-node corresponding to a file identifier.

The following entry points are specified by the Virtual File System interface for performing operations on v-node structures:

Entry Point	Description
<code>vnode_access</code>	Tests a user's permission to access a file.
<code>vnode_close</code>	Releases the resources associated with a v-node.
<code>vnode_create</code>	Creates and opens a new file.
<code>vnode_create_attr</code>	Creates and opens a new file with initial attributes.
<code>vnode_fclear</code>	Releases portions of a file (by zeroing bytes).
<code>vnode_fid</code>	Builds a file identifier for a v-node.
<code>vnode_finfo</code>	Returns pathconf information about a file or file system.
<code>vnode_fsync</code>	Flushes in-memory information and data to permanent storage.
<code>vnode_fsync_range</code>	Flushes in-memory information and data for a given range to permanent storage.
<code>vnode_ftruncate</code>	Decreases the size of a file.
<code>vnode_getacl</code>	Gets information about access control, by retrieving the access control list.
<code>vnode_getattr</code>	Gets the attributes of a file.
<code>vnode_getacl</code>	Gets information about access control by retrieving the ACL. Provides an advanced interface when compared to <code>vnode_getacl</code> .
<code>vnode_hold</code>	Assures that a v-node is not destroyed, by incrementing the v-node's use count.
<code>vnode_ioctl</code>	Performs miscellaneous operations on devices.
<code>vnode_link</code>	Creates a new directory entry for a file.
<code>vnode_lockctl</code>	Sets, removes, and queries file locks.
<code>vnode_lookup</code>	Finds an object by name in a directory.
<code>vnode_map</code>	Associates a file with a memory segment.

Entry Point	Description
<code>vnop_map_lloff</code>	Associates a file with a memory segment using 64 bit offset.
<code>vnop_memcntl</code>	Manages physical attachment of a file.
<code>vnop_mkdir</code>	Creates a directory.
<code>vnop_mknod</code>	Creates a file of arbitrary type.
<code>vnop_open</code>	Gets read and/or write access to a file.
<code>vnop_rdwr</code>	Reads or writes a file.
<code>vnop_rdwr_attr</code>	Reads or writes a file and returns attributes.
<code>vnop_readdir</code>	Reads directory entries in standard format.
<code>vnop_readdir_eofp</code>	Reads directories and returns end of file indication.
<code>vnop_readlink</code>	Reads the contents of a symbolic link.
<code>vnop_rele</code>	Releases a reference to a virtual node (v-node).
<code>vnop_remove</code>	Unlinks a file or directory.
<code>vnop_rename</code>	Renames a file or directory.
<code>vnop_revoke</code>	Revokes access to an object.
<code>vnop_rmdir</code>	Removes a directory.
<code>vnop_seek</code>	Moves the current offset in a file.
<code>vnop_select</code>	Polls a v-node for pending I/O.
<code>vnop_setacl</code>	Sets information about access control for a file.
<code>vnop_setattr</code>	Sets attributes of a file.
<code>vnop_setxcl</code>	Sets information about access control for a file. Provides an advanced interface compared to <code>vnop_setacl</code> .
<code>vnop_strategy</code>	Reads or writes blocks of a file.
<code>vnop_symlink</code>	Creates a symbolic link.
<code>vnop_unmap</code>	Destroys a file or memory association.

#### Related information:

Virtual File System Overview

Virtual File System Kernel Extensions Overview

## vfs\_aclxcntl Entry Point

### Purpose

Implements access-control-specific control operations for a file system.

### Syntax

```
int vfs_aclxcntl (vfsp, vp, cmd, uiop, argsize, crp)
```

```
struct vfs    *vfsp;
struct vnode  *vp;
int           cmd;
struct uio    *uiop;
size_t        argsize;
struct ucred  *crp;
```

### Description

The `vfs_aclxcntl` entry point is invoked to perform various ACL-specific control operations on the underlying physical file system. If a file system is implemented to support this interface, it needs to adhere to the various commands and arguments defined for the interface. A file system implementation can define `cmd` parameter values and corresponding control functions that are specific to the file system. The `cmd` parameter for these functions has values defined globally for all the physical file systems. These control operations can be issued with the ACL library interfaces.

### Parameters

Item	Description
<i>vfsp</i>	Points to the file system for which the control operation is to be issued.
<i>vp</i>	Points to the virtual node pointer to the file path of the file system for which the control operation is being requested.
<i>cmd</i>	Specifies which control operation to perform. Has one of the following values:  <b>ACLNTL_GETACLXYPES</b> Returns the various ACL types supported for the file system instance. This area is of the following structure type: <pre>typedef struct _acl_types_list_t {     uint32_t num_entries; // in the buffer to follow     uint32_t pad; // reserved space     acl_type_t entries[MAX_ACL_TYPES]; } // Array of ACL types acl_types_list_t ;</pre> If the buffer space is not enough to accommodate ACL types supported by the physical file system, <b>errno</b> is set to <b>ENOSPC</b> and the necessary size of the buffer is returned in <i>argsize</i> .  <b>ACLNTL_GETACLXTYPEINFO</b> Returns the characteristics information related to an ACL type for the file system instance. This area is of the following structure type: <pre>typedef struct _acl_type_info_t {     acl_type_t acl_type; } // ACL type for which info is needed uint8_t acl_type_info; } // Start of ACL characteristics data _acl_type_info_t ;</pre> <i>acl_type_info</i> is the start byte of the ACL-related characteristics information. ACL characteristics information depends on the ACL type. ACL characteristics for NFS4 ACL type have the following structure: <pre>typedef struct _nfs4_acl_type_info_t {     uint32_t version; } // Version of this structure uint32_t acl_support; } // Support of Access control entry types. nfs4_acl_type_info_t ;</pre> If the buffer space is not enough to accommodate the ACL types supported by the physical file system, <b>errno</b> is set to <b>ENOSPC</b> and the necessary size of the buffer is returned in <i>argsize</i> .
<i>uiop</i>	Identifies data specific to the control operation. If the <i>cmd</i> parameter has a value of <b>ACLNTL_GETACLXYPES</b> , <i>uiop</i> points to a buffer area where the file system stores the supported ACL types. If the <i>cmd</i> parameter has a value of <b>ACLNTL_GETACLXTYPEINFO</b> , <i>uiop</i> points to a buffer area where the file system stores the ACL characteristics information.
<i>argsize</i>	Identifies the length of the data specified by the <i>arg</i> parameter. This buffer is used to return the necessary buffer size, in case the buffer size provided by the user is not enough.
<i>crp</i>	Points to the <b>cred</b> structure. This structure contains data that the file system can use to validate access permission.

## Execution Environment

The **vfs\_aclxctl** entry point can be called from the process environment only.

## Return Values

Upon successful completion, the **vfs\_aclxctl** entry point returns 0. Nonzero return values are returned from the **/usr/include/sys/errno.h** file to indicate failure.

Item	Description
EACCES	The <i>cmd</i> parameter requires a privilege that the current process does not have.
EINVAL	Indicates that the <i>cmd</i> parameter is not a supported control, or the <i>arg</i> parameter is not a valid argument for the command.
ENOSPC	The input buffer was not sufficient for storing the requested information.

### Related information:

Virtual File System Overview

Virtual File System Kernel Extensions Overview

Logical File System Overview

## vfs\_cntl Entry Point

### Purpose

Implements control operations for a file system.

### Syntax

```
int vfs_cntl (vfsp, cmd, arg, argsize, crp)
struct vfs * vfsp;
int cmd;
caddr_t arg;
unsigned long argsize;
struct ucred * crp;
```

### Parameters

Item	Description
<i>vfsp</i>	Points to the file system for which the control operation is to be issued.
<i>cmd</i>	Specifies which control operation to perform.
<i>arg</i>	Identifies data specific to the control operation.
<i>argsize</i>	Identifies the length of the data specified by the <i>arg</i> parameter.
<i>crp</i>	Points to the <b>cred</b> structure. This structure contains data that the file system can use to validate access permission.

### Description

The **vfs\_cntl** entry point is invoked by the logical file system to request various control operations on the underlying file system. A file system implementation can define file system-specific *cmd* parameter values and corresponding control functions. The *cmd* parameter for these functions should have a minimum value of 32768. These control operations can be issued with the **fscntl** subroutine.

**Note:** The only system-supported control operation is **FS\_EXTENDFS**. This operation increases the file system size and accepts an *arg* parameter that specifies the new size. The **FS\_EXTENDFS** operation ignores the *argsize* parameter.

### Execution Environment

The **vfs\_cntl** entry point can be called from the process environment only.

### Return Values

Item	Description
0	Indicates success.

Non-zero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure. Typical values include:

Item	Description
EINVAL	Indicates that the <code>cmd</code> parameter is not a supported control, or the <code>arg</code> parameter is not a valid argument for the command.
EACCES	Indicates that the <code>cmd</code> parameter requires a privilege that the current process does not have.

#### Related information:

fscntl subroutine

Virtual File System Overview

Virtual File System Kernel Extensions Overview

## vfs\_hold or vfs\_unhold Kernel Service Purpose

Holds or releases a `vfs` structure.

### Syntax

```
#include <sys/vfs.h>
void vfs_hold(vfsp)
struct vfs *vfsp;
int vfs_unhold(vfsp)
struct vfs *vfsp;
```

### Parameter

Item	Description
<i>vfsp</i>	Points to a <code>vfs</code> structure.

### Description

The `vfs_hold` kernel service holds a `vfs` structure and the `vfs_unhold` kernel service releases it. These routines manage a use count for a virtual file system (VFS). A use count greater than 1 prevents the virtual file system from being unmounted.

### Execution Environment

These kernel services can be called from the process environment only.

### Return Values

The `vfs_hold` kernel service has no return value.

The `vfs_unhold` kernel service returns the original value of the hold count.

## vfs\_init Entry Point Purpose

Initializes a virtual file system.

## Syntax

```
int vfs_init ( gfsp)
struct gfs *gfsp;
```

## Parameter

Item	Description
<i>gfsp</i>	Points to a file system's attribute structure.

## Description

The **vfs\_init** entry point is invoked to initialize a file system. It is called when a file system implementation is loaded to perform file system-specific initialization.

The **vfs\_init** entry point is not called through the virtual file system switch. Instead, it is called indirectly by the **gfsadd** kernel service when the **vfs\_init** entry point address is stored in the **gfs** structure passed to the **gfsadd** kernel service as a parameter. (The **vfs\_init** address is placed in the **gfs\_init** field of the **gfs** structure.) The **gfs** structure is defined in the **/usr/include/sys/gfs.h** file.

**Note:** The return value for the **vfs\_init** entry point is passed back as the return value from the **gfsadd** kernel service.

## Execution Environment

The **vfs\_init** entry point can be called from the process environment only.

## Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the **/usr/include/sys/errno.h** file to indicate failure.

### Related reference:

“gfsadd Kernel Service” on page 193

### Related information:

Virtual File System Overview

Virtual File System Kernel Extensions Overview

## vfs\_mount Entry Point

### Purpose

Mounts a virtual file system.

## Syntax

```
int vfs_mount ( vfsp)
struct vfs *vfsp;
struct ucred *crp;
```

## Parameter

Item	Description
<i>vfsp</i>	Points to the newly created <b>vfs</b> structure.
<i>crp</i>	Points to the <b>cred</b> structure. This structure contains data that the file system can use to validate access permission.

## Description

The **vfs\_mount** entry point is called by the logical file system to mount a new file system. This entry point is called after the **vfs** structure is allocated and initialized. Before this structure is passed to the **vfs\_mount** entry point, the logical file system:

- Guarantees the syntax of the **vmount** or **mount** subroutines.
- Allocates the **vfs** structure.
- Resolves the stub to a virtual node (v-node). This is the **vfs\_mntdover** field in the **vfs** structure.
- Initializes the following virtual file system fields:

Field	Description
<b>vfs_flags</b>	<p>Initialized depending on the type of mount. This field takes the following values:</p> <p><b>VFS_MOUNTOK</b> The user has write permission in the stub's parent directory and is the owner of the stub.</p> <p><b>VFS_SUSER</b> The user has root user authority.</p> <p><b>VFS_NOSUID</b> Execution of <b>setuid</b> and <b>setgid</b> programs from this mount are not allowed.</p> <p><b>VFS_NODEV</b> Opens of devices from this mount are not allowed.</p>
<b>vfs_type</b>	Initialized to the / (root) file system type when the <b>mount</b> subroutine is used. If the <b>vmount</b> subroutine is used, the <b>vfs_type</b> field is set to the <i>type</i> parameter supplied by the user. The logical file system verifies the existence of the <i>type</i> parameter.
<b>vfs_ops</b>	Initialized according to the <b>vfs_type</b> field.
<b>vfs_mntdover</b>	Identifies the v-node that refers to the stub path argument. This argument is supplied by the <b>mount</b> or <b>vmount</b> subroutine.
<b>vfs_date</b>	Holds the time stamp. The time stamp specifies the time to initialize the virtual file system.
<b>vfs_number</b>	Indicates the unique number sequence representing this virtual file system.
<b>vfs_mdata</b>	Initialized with the <b>vmount</b> structure supplied by the user. The virtual file system data is detailed in the <b>/usr/include/sys/vmount.h</b> file. All arguments indicated by this field are copied to kernel space.

## Execution Environment

The **vfs\_mount** entry point can be called from the process environment only.

## Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the **/usr/include/sys/errno.h** file to indicate failure.

### Related information:

mount subroutine

Virtual File System Overview

Logical File System Overview

## vfs\_root Entry Point

### Purpose

Returns the root v-node of a virtual file system (VFS).

### Syntax

```
int vfs_root ( vfsp, vpp, crp)
struct vfs *vfsp;
struct vnode **vpp;
struct ucred *crp;
```

### Parameters

Item	Description
<i>vfsp</i>	Points to the <b>vfs</b> structure.
<i>vpp</i>	Points to the place to return the v-node pointer.
<i>crp</i>	Points to the <b>cred</b> structure. This structure contains data that the file system can use to validate access permission.

### Description

The **vfs\_root** entry point is invoked by the logical file system to get a pointer to the root v-node of the file system. When successful, the *vpp* parameter points to the root virtual node (v-node) and the v-node hold count is incremented.

### Execution Environment

The **vfs\_root** entry point can be called from the process environment only.

### Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure.

#### Related information:

Virtual File System Overview

Understanding Data Structures and Header Files for Virtual File Systems

Logical File System Overview

## vfs\_search Kernel Service

### Purpose

Searches the **vfs** list.

### Syntax

```
int vfs_search ( vfs_srchfcn, srchargs)
(int (*vfs_srchfcn)(struct vfs *, caddr_t);
caddr_t srchargs;
```



## Parameters

Item	Description
<i>vfs_srchfcn</i>	Points to a search function. The search function is identified by the <i>vfs_srchfcn</i> parameter. This function is used to examine or modify an entry in the vfs list. The search function is called once for each currently active VFS. If the search function returns a value of 0, iteration through the vfs list continues to the next entry. If the return value is nonzero, <b>vfs_search</b> kernel service returns to its caller, passing back the return value from the search function. When the system invokes this function, the system passes it a pointer to a virtual file system (VFS) and the <i>srchargs</i> parameter.
<i>srchargs</i>	Points to data to be used by the search function. This pointer is not used by the <b>vfs_search</b> kernel service but is passed to the search function.

## Description

The **vfs\_search** kernel service searches the vfs list. This kernel service allows a process outside the file system to search the vfs list. The **vfs\_search** kernel service locks out all activity in the vfs list during a search. Then, the kernel service iterates through the vfs list and calls the search function on each entry.

The search function must not request locks that could result in deadlock. In particular, any attempt to do lock operations on the vfs list or on other VFS structures could produce deadlock.

The performance of the **vfs\_search** kernel service may not be acceptable for functions requiring quick response. Iterating through the vfs list and making an indirect function call for each structure is inherently slow.

## Execution Environment

The **vfs\_search** kernel service can be called from the process environment only.

## Return Values

This kernel service returns the value returned by the last call to the search function.

## Related Information

### **vfs\_statfs** Entry Point

#### Purpose

Returns virtual file system statistics.

#### Syntax

```
int vfs_statfs ( vfsp, stafsp, crp )
struct vfs *vfsp;
struct statfs *stafsp;
struct ucred *crp;
```

#### Parameters

Item	Description
<i>vfsp</i>	Points to the <b>vfs</b> structure being queried. This structure is defined in the <code>/usr/include/sys/vfs.h</code> file.
<i>stafsp</i>	Points to a <b>statfs</b> structure. This structure is defined in the <code>/usr/include/sys/statfs.h</code> file.
<i>crp</i>	Points to the <b>cred</b> structure. This structure contains data that the file system can use to validate access permission.

## Description

The **vfs\_statfs** entry point is called by the logical file system to obtain file system characteristics. Upon return, the **vfs\_statfs** entry point has filled in the following fields of the **statfs** structure:

Field	Description
<code>f_blocks</code>	Specifies the number of blocks.
<code>f_files</code>	Specifies the total number of file system objects.
<code>f_bsize</code>	Specifies the file system block size.
<code>f_bfree</code>	Specifies the number of free blocks.
<code>f_ffree</code>	Specifies the number of free file system objects.
<code>f_fname</code>	Specifies a 32-byte string indicating the file system name.
<code>f_fpack</code>	Specifies a 32-byte string indicating a pack ID.
<code>f_name_max</code>	Specifies the maximum length of an object name.

Fields for which a **vfs** structure has no values are set to 0.

## Execution Environment

The **vfs\_statfs** entry point can be called from the process environment only.

## Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure.

### Related information:

[statfs subroutine](#)

[Virtual File System Overview](#)

[Virtual File System Kernel Extensions Overview](#)

## vfs\_sync Entry Point

### Purpose

Requests that file system changes be written to permanent storage.

### Syntax

```
int vfs_sync (* gfsp)
struct gfs *gfsp;
```

### Parameter

Item	Description
<i>gfsp</i>	Points to a <b>gfs</b> structure. The <b>gfs</b> structure describes the file system type. This structure is defined in the <code>/usr/include/sys/gfs.h</code> file.

## Description

The **vfs\_sync** entry point is used by the logical file system to force all data associated with a particular virtual file system type to be written to its storage. This entry point is used to establish a known consistent state of the data.

**Note:** The **vfs\_sync** entry point is called once per file system type rather than once per virtual file system.

## Execution Environment

The **vfs\_sync** entry point can be called from the process environment only.

## Return Values

The **vfs\_sync** entry point is advisory. It has no return values.

### Related information:

sync subroutine

Virtual File System Overview

Virtual File System Kernel Extensions Overview

Logical File System Overview

## vfs\_umount Entry Point

### Purpose

Unmounts a virtual file system.

### Syntax

```
int vfs_umount ( vfsp, crp )
struct vfs *vfsp;
struct ucred *crp;
```

### Parameters

Item	Description
<i>vfsp</i>	Points to the <b>vfs</b> structure being unmounted. This structure is defined in the <code>/usr/include/sys/vfs.h</code> file.
<i>crp</i>	Points to the <b>cred</b> structure. This structure contains data that the file system can use to validate access permission.

## Description

The **vfs\_umount** entry point is called to unmount a virtual file system. The logical file system performs services independent of the virtual file system that initiate the unmounting. The logical file system services:

- Guarantee the syntax of the **umount** subroutine.
- Perform permission checks:
  - If the *vfsp* parameter refers to a device mount, then the user must have root user authority to perform the operation.

- If the *vfsp* parameter does not refer to a device mount, then the user must have root user authority or write permission in the parent directory of the mounted-over virtual node (v-node), as well as write permission to the file represented by the mounted-over v-node.
- Ensure that the virtual file system being unmounted contains no mount points for other virtual file systems.
- Ensure that the root v-node is not in use except for the mount. The root v-node is also referred to as the mounted v-node.
- Clear the *v\_mvfsp* field in the stub v-node. This prevents lookup operations already in progress from traversing the soon-to-be unmounted mount point.

The logical file system assumes that, if necessary, successful **vfs\_umount** entry point calls free the root v-node. An error return from the **vfs\_umount** entry point causes the mount point to be re-established. A 0 (zero) returned from the **vfs\_umount** entry point indicates the routine was successful and that the **vfs** structure was released.

## Execution Environment

The **vfs\_umount** entry point can be called from the process environment only.

## Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure.

### Related information:

umount subroutine

vmount subroutine

Understanding Data Structures and Header Files for Virtual File Systems

## vfs\_vget Entry Point

### Purpose

Converts a file identifier into a virtual node (v-node).

### Syntax

```
int vfs_vget ( vfsp, vpp, fidp, crp)
struct vfs *vfsp;
struct vnode **vpp;
struct fileid *fidp;
struct ucred *crp;
```

### Parameters

Item	Description
<i>vfsp</i>	Points to the virtual file system that is to contain the v-node. Any returned v-node should belong to this virtual file system.
<i>vpp</i>	Points to the place to return the v-node pointer. This is set to point to the new v-node. The fields in this v-node should be set as follows: <ul style="list-style-type: none"> <li><b>v_vntype</b> The type of v-node dependent on private data.</li> <li><b>v_count</b> Set to at least 1 (one).</li> <li><b>v_pdata</b> If a new file, set to the private data for this file system.</li> </ul>
<i>fidp</i>	Points to a file identifier. This is a file system-specific file identifier that must conform to the <b>fileid</b> structure. <p><b>Note:</b> If the <i>fidp</i> parameter is invalid, the <i>vpp</i> parameter should be set to a null value by the <b>vfs_vget</b> entry point.</p>
<i>crp</i>	Points to the <b>cred</b> structure. This structure contains data that the file system can use to validate access permission.

## Description

The **vfs\_vget** entry point is called to convert a file identifier into a v-node. This entry point uses information in the *vfsp* and *fidp* parameters to create a v-node or attach to an existing v-node. This v-node represents, logically, the same file system object as the file identified by the *fidp* parameter.

If the v-node already exists, successful operation of this entry point increments the v-node use count and returns a pointer to the v-node. If the v-node does not exist, the **vfs\_vget** entry point creates it using the **vn\_get** kernel service and returns a pointer to the new v-node.

## Execution Environment

The **vfs\_vget** entry point can be called from the process environment only.

## Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the **/usr/include/sys/errno.h** file to indicate failure. A typical value includes:

Item	Description
EINVAL	Indicates that the remote virtual file system specified by the <i>vfsp</i> parameter does not support chained mounts.

### Related reference:

“vn\_get Kernel Service” on page 580

### Related information:

access subroutine

Virtual File System Overview

## vnop\_access Entry Point

### Purpose

Requests validation of user access to a virtual node (v-node).

### Syntax

```
int vnop_access ( vp, mode, who, crp)
struct vnode *vp;
```

```
int mode;
int who;
struct ucred *crp;
```

## Parameters

Item	Description
<i>vp</i>	Points to the v-node.
<i>mode</i>	Identifies the access mode.
<i>who</i>	Specifies the IDs for which to check access. This parameter should be one of the following values, which are defined in the <code>/usr/include/sys/access.h</code> file: <ul style="list-style-type: none"> <li><b>ACC_SELF</b> Determines if access is permitted for the current process. The effective user and group IDs and the supplementary group ID of the current process are used for the calculation.</li> <li><b>ACC_ANY</b> Determines if the specified access is permitted for any user, including the object owner. The <i>mode</i> parameter must contain only one of the valid modes.</li> <li><b>ACC_OTHERS</b> Determines if the specified access is permitted for any user, excluding the owner. The <i>mode</i> parameter must contain only one of the valid modes.</li> <li><b>ACC_ALL</b> Determines if the specified access is permitted for all users. (This is a useful check to make when files are to be written blindly across networks.) The <i>mode</i> parameter must contain only one of the valid modes.</li> </ul>
<i>crp</i>	Points to the <code>cred</code> structure. This structure contains data that the file system can use to validate access permission.

## Description

The `vnop_access` entry point is used by the logical volume file system to validate access to a v-node. This entry point is used to implement the `access` subroutine. The v-node is held for the duration of the `vnop_access` entry point. The v-node count is unchanged by this entry point.

In addition, the `vnop_access` entry point is used for permissions checks from within the file system implementation. The valid types of access are listed in the `/usr/include/sys/access.h` file. Current modes are read, write, execute, and existence check.

**Note:** The `vnop_access` entry point must ensure that write access is not requested on a read-only file system.

## Execution Environment

The `vnop_access` entry point can be called from the process environment only.

## Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure. A typical value includes:

Item	Description
EACCES	Indicates no access is allowed.

#### Related information:

access subroutine

Virtual File System Overview

Virtual File System Kernel Extensions Overview

## vnop\_close Entry Point

### Purpose

Closes a file associated with a v-node (virtual node).

### Syntax

```
int vnop_close ( vp, flag, vinfo, crp)
struct vnode *vp;
int flag;
caddr_t vinfo;
struct ucred *crp;
```

### Parameters

Item	Description
<i>vp</i>	Points to the v-node.
<i>flag</i>	Identifies the flag word from the file pointer.
<i>vinfo</i>	This parameter is not used.
<i>crp</i>	Points to the <b>cred</b> structure. This structure contains data that the file system can use to validate access permission.

### Description

The **vnop\_close** entry point is used by the logical file system to announce that the file associated with a given v-node is now closed. The v-node continues to remain active but will no longer receive read or write requests through the **vnop\_rdwr** entry point.

A **vnop\_close** entry point is called only when the use count of an associated file structure entry goes to 0 (zero).

**Note:** The v-node is held over the duration of the **vnop\_close** entry point.

### Execution Environment

The **vnop\_close** entry point can be called from the process environment only.

### Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure.

**Note:** The `vnop_close` entry point may fail and an error will be returned to the application. However, the v-node is considered closed.

**Related reference:**

“`vnop_open` Entry Point” on page 675

**Related information:**

close subroutine

Virtual File System Overview

## vnop\_create Entry Point

### Purpose

Creates a new file.

### Syntax

```
int vnop_create (dp, vpp, flag, pname, mode, vinfop, crp)
struct vnode * dp;
struct vnode ** vpp;
int flag;
char * pname;
int mode;
caddr_t * vinfop;
struct ucred * crp;
```

### Parameters

Item	Description
<i>dp</i>	Points to the virtual node (v-node) of the parent directory.
<i>vpp</i>	Points to the place in which the pointer to a v-node for the newly created file is returned.
<i>flag</i>	Specifies an integer flag word. The <code>vnop_create</code> entry point uses this parameter to open the file.
<i>pname</i>	Points to the name of the new file.
<i>mode</i>	Specifies the mode for the new file.
<i>vinfop</i>	This parameter is unused.
<i>crp</i>	Points to the <code>cred</code> structure. This structure contains data that the file system can use to validate access permission.

### Description

The `vnop_create` entry point is invoked by the logical file system to create a regular (v-node type `VREG`) file in the directory specified by the `dp` parameter. (Other v-node operations create directories and special files.) Virtual node types are defined in the `/usr/include/sys/vnode.h` file. The v-node of the parent directory is held during the processing of the `vnop_create` entry point.

To create a file, the `vnop_create` entry point does the following:

- Opens the newly created file.
- Checks that the file system associated with the directory is not read-only.

**Note:** The logical file system calls the `vnop_lookup` entry point before calling the `vnop_create` entry point.



## Execution Environment

The `vnop_create` entry point can be called from the process environment only.

## Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure.

### Related reference:

“`vnop_lookup` Entry Point” on page 670

### Related information:

Virtual File System Overview

Virtual File System Kernel Extensions Overview

## `vnop_create_attr` Entry Point

### Purpose

Creates a new file.

### Syntax

```
int vnop_create_attr (dvp, vpp, flags, name, vap, vcf, finfop, crp) struct vnode *dvp; struct vnode *vpp; int flags; char *name; struct vattr *vap; int vcf; caddr_t finfop; struct ucred *crp;
```

### Parameters

Item	Description
<i>dvp</i>	Points to the directory vnode.
<i>vpp</i>	Points to the newly created vnode pointer.
<i>flags</i>	Specifies file creation flags.
<i>name</i>	Specifies the name of the file to create.
<i>vattr</i>	Points to the initial attributes.
<i>vcf</i>	Specifies create flags.
<i>finfop</i>	Specifies address of finfo field.
<i>crp</i>	Specifies user's credentials.

### Description

The `vnop_create_attr` entry point is used to create a new file. This operation is similar to the `vnop_create` entry point except that the initial file attributes are passed in a `vattr` structure.

The `va_mask` field in the `vattr` structure identifies which attributes are to be applied. For example, if the `AT_SIZE` bit is set, then the file system should use `va_size` for the initial file size. For all `vnop_create_attr` calls, at least `AT_TYPE` and `AT_MODE` must be set.

The `vcf` parameter controls how the new vnode is to be activated. If `vcf` is set to `VC_OPEN`, then the new object should be opened. If `vcf` is `VC_LOOKUP`, then the new object should be created, but not opened. If `vcf` is `VC_DEFAULT`, then the new object should be created, but the vnode for the object is not activated.

File systems that do not define `GFS_VERSION421` in their `gfs` flags do not need to supply a `vnop_create_attr` entry point. The logical file system will funnel all creation requests through the old `vnop_create` entry point.

## Execution Environment

The `vnop_create_attr` entry point can be called from the process environment only.

## Return Values

Item	Description
Zero	Indicates a successful operation; <i>*vpp</i> contains a pointer to the new vnode.
Nonzero	Indicates that the operation failed; return values should be chosen from the <code>/usr/include/sys/errno.h</code> file.

## vnop\_fclear Entry Point

### Purpose

Releases portions of a file.

### Syntax

```
int vnop_fclear (vp, flags, offset, len, vinfo, crp)
struct vnode * vp;
int flags;
offset_t offset;
offset_t len;
caddr_t vinfo;
struct ucred * crp;
```

### Parameters

Item	Description
<i>vp</i>	Points to the virtual node (v-node) of the file.
<i>flags</i>	Identifies the flags from the open file structure.
<i>offset</i>	Indicates where to start clearing in the file.
<i>len</i>	Specifies the length of the area to be cleared.
<i>vinfo</i>	This parameter is unused.
<i>crp</i>	Points to the <code>cred</code> structure. This structure contains data that the file system can use to validate access permission.

### Description

The `vnop_fclear` entry point is called from the logical file system to clear bytes in a file, returning whole free blocks to the underlying file system. This entry point performs the clear regardless of whether the file is mapped.

Upon completion of the `vnop_fclear` entry point, the logical file system updates the file offset to reflect the number of bytes cleared. Also upon completion, if either the starting or ending offset is past the starting end of file, the file is extended.

## Execution Environment

The `vnop_fclear` entry point can be called from the process environment only.

## Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure.

**Related information:**

`fclear` subroutine

Virtual File System Overview

## **vnop\_fid Entry Point**

### **Purpose**

Builds a file identifier for a virtual node (v-node).

### **Syntax**

```
int vnop_fid ( vp, fidp, crp)
struct vnode *vp;
struct fileid *fidp;
struct ucred *crp;
```

### **Parameters**

Item	Description
<i>vp</i>	Points to the v-node that requires the file identifier.
<i>fidp</i>	Points to where to return the file identifier.
<i>crp</i>	Points to the <code>cred</code> structure. This structure contains data that the file system can use to validate access permission.

### **Description**

The `vnop_fid` entry point is invoked to build a file identifier for the given v-node. This file identifier must contain sufficient information to find a v-node that represents the same file when it is presented to the `vfs_get` entry point.

### **Execution Environment**

The `vnop_fid` entry point can be called from the process environment only.

### **Return Values**

Item	Description
0	Indicates success.

Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure.

**Related information:**

Virtual File System Overview

Virtual File System Kernel Extensions Overview

Logical File System Overview

## **vnop\_finfo Entry Point**

### **Purpose**

Returns information about a file.

## Syntax

```
int vnop_finfo (vp, cmd, bufp, length, crp) struct vnode *vp; int cmd; void *bufp; int length; struct ucred *crp;
```

## Parameters

Item	Description
<i>vp</i>	Points to the vnode to be queried.
<i>cmd</i>	Specifies the command parameter.
<i>bufp</i>	Points to the buffer for the information.
<i>length</i>	Specifies the length of the buffer.
<i>crp</i>	Specifies user's credentials.

## Description

The **vnop\_finfo** entry point is used to query a file system. It is used primarily to implement the **pathconf** and **fpathconf** subroutines. The **command** parameter defines what type of query is being done. The query commands and the associated data structures are defined in **<sys/finfo.h>**. If the file system does not support the particular query, it should return ENOSYS.

File systems that do not define GFS\_VERSION421 in their gfs flags do not need to supply a **vnop\_finfo** entry point. If the command is FI\_PATHCONF, then the logical file system returns generic pathconf information. If the query is other than FI\_PATHCONF, then the request fails with EINVAL.

## Execution Environment

The **vnop\_finfo** entry point can be called from the process environment only.

## Return Values

Item	Description
Zero	Indicates a successful operation.
Nonzero	Indicates that the operation failed; return values should be chosen from the <b>/usr/include/sys/errno.h</b> file.

### Related information:

pathconf, fpathconf

Virtual File System Overview

Logical File System Overview

## vnop\_fsync, vnop\_fsync\_range Entry Points

### Purpose

Flushes file data from memory to disk.

## Syntax

```
int vnop_fsync ( vp, flags, vinfo, crp)
struct vnode *vp;
long flags;
long vinfo;
struct ucred *crp;
```

```
int vnop_fsync_range ( vp, flags, vinfo, offset, length, crp)
struct vnode *vp;
```

```

long flags;
long vinfo;
offset_t offset;
offset_t length;
struct ucred *crp;

```

## Parameters

Item	Description
<i>vp</i>	Points to the virtual node (v-node) of the file.
<i>flags</i>	Identifies flags from the open file and the flags that govern the action to be taken. It can be one of the following values: <p><b>FDATASYNC</b> The changed data in the range specified by the <i>offset</i> and <i>length</i> parameters is written to the storage. If the metadata of the file is changed and this changed metadata must read the data, the metadata is also written to the storage. Otherwise, the metadata is not updated.</p> <p><b>FFILESYNC</b> The changed data in the range specified by the <i>offset</i> and <i>length</i> parameters is written to the storage. If any metadata is changed, all of the changed user data is written to the storage. Metadata changes and file attributes including time stamps are also written to the storage.</p> <p><b>FNOCACHE</b> The changed data is written to the storage similar to the <b>FDATASYNC</b> flag value. The full pages in the range specified by the <i>offset</i> and <i>length</i> parameters are removed from the memory cache. The pages are removed from the cache even if the pages are not changed. This operation is also applicable to the files that are open only for reading.</p>
<i>vinfo</i>	This parameter is currently not used.
<i>offset</i>	Specifies the starting offset in the file of the data to be flushed.
<i>length</i>	Specifies the length of the data to be flushed. If you specify the value as zero, all cached data is flushed.
<i>crp</i>	Points to the <b>cred</b> structure. This structure contains data that the file system can use to validate access permission.

## Description

The **vnop\_fsync** entry point is called by the logical file system to request that all modifications associated with a given v-node to be flushed out to permanent storage. This must be done synchronously so that the caller can assure that all I/O has completed successfully. The **vnop\_fsync\_range** entry point provides the same function but limits the data to be written to a specified range in the file.

## Execution Environment

The **vnop\_fsync** and **vnop\_fsync\_range** entry points can be called from the process environment only.

## Return Values

Item	Description
0	Indicates success.

Nonzero values are returned from the **/usr/include/sys/errno.h** file to indicate failure.

### Related information:

[fsync subroutine](#)

[Virtual File System Kernel Extensions Overview](#)

[Logical File System Overview](#)

## vnop\_ftrunc Entry Point

### Purpose

Truncates a file.

### Syntax

```
int vnop_ftrunc (vp, flags, length, vinfo, crp)
struct vnode * vp;
int flags;
offset_t length;
caddr_t vinfo;
struct ucred * crp;
```

### Parameters

Item	Description
<i>vp</i>	Points to the virtual node (v-node) of the file.
<i>flags</i>	Identifies flags from the open file structure.
<i>length</i>	Specifies the length to which the file should be truncated.
<i>vinfo</i>	This parameter is unused.
<i>crp</i>	Points to the <b>cred</b> structure. This structure contains data that the file system can use to validate access permission.

### Description

The **vnop\_ftrunc** entry point is invoked by the logical file system to decrease the length of a file by truncating it. This operation is unsuccessful if any process other than the caller has locked a portion of the file past the specified offset.

### Execution Environment

The **vnop\_ftrunc** entry point can be called from the process environment only.

### Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the **/usr/include/sys/errno.h** file to indicate failure.

#### Related information:

ftruncate subroutine

Virtual File System Overview

Logical File System Overview

## vnop\_getacl Entry Point

### Purpose

Retrieves the access control list (ACL) for a file.

### Syntax

```
#include <sys/acl.h>
```

```
int vnop_getacl ( vp, uiop, crp)
struct vnode *vp;
struct uio *uiop;
struct ucred *crp;
```

## Description

The `vnop_getacl` entry point is used by the logical file system to retrieve the access control list (ACL) for a file to implement the `getacl` subroutine.

## Parameters

Item	Description
<i>vp</i>	Specifies the virtual node (v-node) of the file system object.
<i>uiop</i>	Specifies the <code>uio</code> structure that defines the storage for the ACL.
<i>crp</i>	Points to the <code>cred</code> structure. This structure contains data that the file system can use to validate access permission.

## Execution Environment

The `vnop_getacl` entry point can be called from the process environment only.

## Return Values

Item	Description
0	Indicates a successful operation.

Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure. A valid value includes:

Item	Description
ENOSPC	Indicates that the buffer size specified in the <code>uiop</code> parameter was not large enough to hold the ACL. If this is the case, the first word of the user buffer (data in the <code>uio</code> structure specified by the <code>uiop</code> parameter) is set to the appropriate size.

### Related reference:

“uio Structure” on page 639

### Related information:

`chacl` subroutine

`statacl` subroutine

Virtual File System Overview

## vnop\_getattr Entry Point Purpose

Gets the attributes of a file.

## Syntax

```
int vnop_getattr ( vp, vap, crp)
struct vnode *vp;
struct vattr *vap;
struct ucred *crp;
```

## Parameters

Item	Description
<i>vp</i>	Specifies the virtual node (v-node) of the file system object.
<i>vap</i>	Points to a <b>vattr</b> structure.
<i>crp</i>	Points to the <b>cred</b> structure. This structure contains data that the file system can use to validate access permission.

## Description

The **vnop\_getattr** entry point is called by the logical file system to retrieve information about a file. The **vattr** structure indicated by the *vap* parameter contains all the relevant attributes of the file. The **vattr** structure is defined in the `/usr/include/sys/vattr.h` file. This entry point is used to implement the **stat**, **fstat**, and **lstat** subroutines.

**Note:** The indicated v-node is held for the duration of the **vnop\_getattr** subroutine.

## Execution Environment

The **vnop\_getattr** entry point can be called from the process environment only.

## Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure.

### Related information:

statx subroutine

Virtual File System Overview

Virtual File System Kernel Extensions Overview

## vnop\_getxacl Entry Point

### Purpose

Retrieves the access control list (ACL) for a file. This is an advanced version of **vnop\_getacl** interface.

### Syntax

```
#include <sys/acl.h>
int vnop_getxacl (vp, ctl_flags, acl_type, uiop, acl_len, mode_info, crp)
```

```
struct vnode    *vp;
uint64_t       ctl_flags;
acl_type_t     *acl_type;
struct uiop    *uiop;
size_t         *acl_len;
mode_t         *mode_info;
struct ucred   *crp;
```

### Description

The **vnop\_getxacl** entry point retrieves the access control list (ACL) for a file system object. It is an advanced version of **vnop\_getacl** interface and provides for ACL-type-based operations. Note that this interface can be used to obtain the ACL type and length information, without actually retrieving the ACL data (see the *ctl\_flags* description for more details).



## Parameters

Item	Description
<i>vp</i>	Specifies the virtual node (v-node) of the file system object.
<i>acl_type</i>	Points to buffer space for file systems to return the ACL type associated with the file system object. The value should normally be set to <b>ACL_ANY</b> or 0 when the call is made. Some physical file systems can solicit ACL requests for a particular ACL type. In such cases, the caller provides the ACL type requested in this buffer. <b>Note:</b> The latter issue is file system implementation specific. For example, when ACL information is requested with an input ACL type, a physical file system might return an error if the existing ACL associated with the file system object is of a different ACL type. Or, the file system might emulate an ACL of the type requested and return.
<i>acl_len</i>	Pointer to a <i>length</i> variable. The space pointed to is used as an input, as well as output, parameter. As input, the value will indicate the size of buffer <i>uiop</i> . When the call returns, this space holds the actual length of the ACL (true for when the call is successful or when the call fails with <b>errno</b> set to <b>ENOSPC</b> ).
<i>ctl_flags</i>	A 64-bit bit mask that provides control over the ACL retrieval and for any future variations in the interface. The following value is defined for these flags: <b>GET_ACLINFO_ONLY</b> Gets only the ACL type and length information from the underlying file system. When this bit is set, arguments such as <i>mode_info</i> can be set to NULL. All other cases must be valid buffer pointers or else an error is returned. If this bit is not specified, all the other information about the ACL (such as ACL data and mode information) is returned.
<i>uiop</i>	Specifies the <b>uio</b> structure that provides space for the store of the ACL.
<i>mode_info</i>	This value indicates any mode word information that needs to be retrieved for the file system object as part of this ACL get operation.
<i>crp</i>	Points to the <b>cred</b> structure. This structure contains data that the file system can use to validate access permission.

## Execution Environment

The **vnop\_getxacl** entry point can be called from the process environment only.

## Return Values

Upon successful completion, the **vnop\_getxacl** entry point returns 0. Nonzero return values are returned from the **/usr/include/sys/errno.h** file to indicate failure.

Item	Description
<b>ENOSPC</b>	Indicates that the buffer size specified in the <i>uiop</i> parameter was not large enough to hold the ACL.

**Note:** This list of error numbers is not complete and is dependent on the particular physical file system implementation supporting the ACL.

### Related reference:

“uio Structure” on page 639

### Related information:

chacl subroutine

Virtual File System Overview

## vnop\_hold Entry Point

### Purpose

Assures that a virtual node (v-node) is not destroyed.

## Syntax

```
int vnop_hold ( vp)
struct vnode *vp;
```

## Parameter

Item	Description
<i>vp</i>	Points to the v-node.

## Description

The **vnop\_hold** entry point increments the `v_count` field, the hold count on the v-node, and the v-node's underlying g-node (generic node). This incrementation assures that the v-node is not deallocated.

## Execution Environment

The **vnop\_hold** entry point can be called from the process environment only.

## Return Values

The **vnop\_hold** entry point cannot fail and therefore has no return values.

### Related information:

Virtual File System Overview

## vnop\_ioctl Entry Point

### Purpose

Requests I/O control operations on special files.

## Syntax

```
int vnop_ioctl (vp, cmd, arg, flags, ext, crp)
struct vnode * vp;
int cmd;
caddr_t arg;
int flags, ext;
struct ucred * crp;
```

## Parameters

Item	Description
<i>vp</i>	Points to the virtual node (v-node) on which to perform the operation.
<i>cmd</i>	Identifies the specific command. Common operations for the <b>ioctl</b> subroutine are defined in the <code>/usr/include/sys/ioctl.h</code> file. The file system implementation can define other <b>ioctl</b> operations.
<i>arg</i>	Defines a command-specific argument. This parameter can be a single word or a pointer to an argument (or result structure).
<i>flags</i>	Identifies flags from the open file structure.
<i>ext</i>	Specifies the extended parameter passed by the <b>ioctl</b> subroutine. The <b>ioctl</b> subroutine always sets the <i>ext</i> parameter to 0.
<i>crp</i>	Points to the <b>cred</b> structure. This structure contains data that the file system can use to validate access permission.

## Description

The `vnop_ioctl` entry point is used by the logical file system to perform miscellaneous operations on special files. If the file system supports special files, the information is passed down to the `ddioctl` entry point of the device driver associated with the given v-node.

## Execution Environment

The `vnop_ioctl` entry point can be called from the process environment only.

## Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure. A valid value includes:

Item	Description
EINVAL	Indicates the file system does not support the entry point.

### Related information:

`ioctl` subroutine

Logical File System Overview

## `vnop_link` Entry Point

### Purpose

Requests a hard link to a file.

### Syntax

```
int vnop_link ( vp, dp, name, crp)
struct vnode *vp;
struct vnode *dp;
caddr_t *name;
struct ucred *crp;
```

### Parameters

Item	Description
<i>vp</i>	Points to the virtual node (v-node) to link to. This v-node is held for the duration of the linking process.
<i>dp</i>	Points to the v-node for the directory in which the link is created. This v-node is held for the duration of the linking process.
<i>name</i>	Identifies the new name of the entry.
<i>crp</i>	Points to the <code>cred</code> structure. This structure contains data that the file system can use to validate access permission.

## Description

The `vnop_link` entry point is invoked to create a new hard link to an existing file as part of the link subroutine. The logical file system ensures that the *dp* and *vp* parameters reside in the same virtual file system, which is not read-only.

## Execution Environment

The `vnop_link` entry point can be called from the process environment only.

## Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure.

### Related information:

Virtual File System Overview

Virtual File System Kernel Extensions Overview

Logical File System Overview

## vnop\_lockctl Entry Point

### Purpose

Sets, checks, and queries record locks.

### Syntax

```
int vnop_lockctl (vp, offset, lckdat, cmd, retry_fn, retry_id, crp)
struct vnode * vp;
offset_t offset;
struct eflock * lckdat;
int cmd;
int (* retry_fn)();
caddr_t retry_id;
struct ucred * crp;
```

### Parameters

Item	Description
<i>vp</i>	Points to the file's virtual node (v-node).
<i>offset</i>	Indicates the file offset from the open file structure. This parameter is used to establish where the lock region begins.
<i>lckdat</i>	Points to the <b>eflock</b> structure. This structure describes the lock operation to perform.
<i>cmd</i>	Identifies the type of lock operation the <b>vnop_lockctl</b> entry point is to perform. It is a bit mask that takes the following lock-control values:  <b>SETFLCK</b> If set, performs a lock set or clear. If clear, returns the lock information. The <code>l_type</code> field in the <b>eflock</b> structure indicates whether a lock is set or cleared.  <b>SLPFLCK</b> If the lock is unavailable immediately, wait for it. This is only valid when the <b>SETFLCK</b> flag is set.
<i>retry_fn</i>	Points to a subroutine that is called when a lock is retried. This subroutine is not used if the lock is granted immediately. <b>Note:</b> If the <i>retry_fn</i> parameter is not a null value, the <b>vnop_lockctl</b> entry point will not sleep, regardless of the <b>SLPFLCK</b> flag.
<i>retry_id</i>	Points to the location where a value can be stored. This value can be used to correlate a retry operation with a specific lock or set of locks. The retry value is only used in conjunction with the <i>retry_fn</i> parameter. <b>Note:</b> This value is an opaque value and should not be used by the caller for any purpose other than a lock correlation. (This value should not be used as a pointer.)
<i>crp</i>	Points to the <b>cred</b> structure. This structure contains data that the file system can use to validate access permission.

## Description

The **vnop\_lockctl** entry point is used to request record locking. This entry point uses the information in the **eflock** structure to implement record locking.

If a requested lock is blocked by an existing lock, the **vnop\_lockctl** entry point should establish a sleeping lock with the retry subroutine address (specified by the *retry\_fn* parameter) stored in the entry point. The **vnop\_lockctl** entry point then returns a correlating ID value to the caller (in the *retry\_id* parameter), along with an exit value of **EAGAIN**. When the sleeping lock is later awakened, the retry subroutine is called with the *retry\_id* parameter as its argument.

## eflock Structure

The **eflock** structure is defined in the **/usr/include/sys/flock.h** file and includes the following fields:

Field	Description
<b>l_type</b>	Specifies type of lock. This field takes the following values: <b>F_RDLCK</b> Indicates read lock. <b>F_WRLCK</b> Indicates write lock. <b>F_UNLCK</b> Indicates unlock this record. A value of <b>F_UNLCK</b> starting at 0 until 0 for a length of 0 means unlock all locks on this file. Unlocking is done automatically when a file is closed.
<b>l_whence</b>	Specifies location that the <b>l_start</b> field offsets.
<b>l_start</b>	Specifies offset from the <b>l_whence</b> field.
<b>l_len</b>	Specifies length of record. If this field is 0, the remainder of the file is specified.
<b>l_vfs</b>	Specifies virtual file system that contains the file.
<b>l_sysid</b>	Specifies value that uniquely identifies the host for a given virtual file system. This field must be filled in before the call to the <b>vnop_lockctl</b> entry point.
<b>l_pid</b>	Specifies process ID (PID) of the lock owner. This field must be filled in before the call to the <b>vnop_lockctl</b> entry point.

## Execution Environment

The **vnop\_lockctl** entry point can be called from the process environment only.

## Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the **/usr/include/sys/errno.h** file to indicate failure. Valid values include:

Item	Description
<b>EAGAIN</b>	Indicates a blocking lock exists and the caller did not use the <b>SLPFLCK</b> flag to request that the operation sleep.
<b>ERRNO</b>	Returns an error number from the <b>/usr/include/sys/errno.h</b> file on failure.

## Related information:

Virtual File System Overview

Logical File System Overview

## vnop\_lookup Entry Point

### Purpose

Returns a v-node for a given name in a directory.

### Syntax

```
int vnop_lookup (dvp, vpp, name, vattrp , crp)
struct vnode * dvp;
struct vnode ** vpp;
char * name;
struct vattr * vattrp;
struct ucred * crp;
```

### Parameters

Item	Description
<i>dvp</i>	Points to the virtual node (v-node) of the directory to be searched. The logical file system verifies that this v-node is of a VDIR type.
<i>name</i>	Points to a null-terminated character string containing the file name to look up.
<i>vattrp</i>	Points to a <b>vattr</b> structure. If this pointer is NULL, no action is required of the file system implementation. If it is not NULL, the attributes of the file specified by the <i>name</i> parameter are returned at the address passed in the <i>vattrp</i> parameter.
<i>vpp</i>	Points to the place to which to return the v-node pointer, if the pointer is found. Otherwise, a null character should be placed in this memory location.
<i>crp</i>	Points to the <b>cred</b> structure. This structure contains data that the file system can use to validate access permission.

### Description

The **vnop\_lookup** entry point is invoked by the logical file system to find a v-node. It is used by the kernel to convert application-given path names to the v-nodes that represent them.

The use count in the v-node specified by the *dvp* parameter is incremented for this operation, and it is not decremented by the file system implementation.

If the name is found, a pointer to the desired v-node is placed in the memory location specified by the *vpp* parameter, and the v-node hold count is incremented. (In this case, this entry point returns 0.) If the file name is not found, a null character is placed in the *vpp* parameter, and the function returns a **ENOENT** value. Errors are reported with a return code from the **/usr/include/sys/errno.h** file. Possible errors are usually specific to the particular virtual file system involved.

### Execution Environment

The **vnop\_lookup** entry point can be called from the process environment only.

### Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure.

#### Related information:

Virtual File System Overview

Virtual File System Kernel Extensions Overview

Logical File System Overview

## vnop\_map Entry Point

### Purpose

Validates file mapping requests.

### Syntax

```
int vnop_map (vp, addr, length, offset, flags, crp)
struct vnode * vp;
caddr_t addr;
uint length;
uint offset;
uint flags;
struct ucred * crp;
```

### Parameters

**Note:** The *addr*, *offset*, and *length* parameters are unused in the current implementation. The file system is expected to store the segment ID with the file in the `gn_seg` field of the g-node for the file.

Item	Description
<i>vp</i>	Points to the virtual node (v-node) of the file.
<i>addr</i>	Identifies the location within the process address space where the mapping is to begin.
<i>length</i>	Specifies the maximum size to be mapped.
<i>offset</i>	Specifies the location within the file where the mapping is to begin.
<i>flags</i>	Identifies what type of mapping to perform. This value is composed of bit values defined in the <code>/usr/include/sys/shm.h</code> file. The following values are of particular interest to file system implementations:  <b>SHM_RDONLY</b> The virtual memory object is read-only.  <b>SHM_COPY</b> The virtual memory object is copy-on-write. If this value is set, updates to the segment are deferred until an <b>fsync</b> operation is performed on the file. If the file is closed without an <b>fsync</b> operation, the modifications are discarded. The application that called the <b>vnop_map</b> entry point is also responsible for calling the <b>vnop_fsync</b> entry point. <b>Note:</b> Mapped segments do not reflect modifications made to a copy-on-write segment.
<i>crp</i>	Points to the <b>cred</b> structure. This structure contains data that applications can use to validate access permission.

### Description

The **vnop\_map** entry point is called by the logical file system to validate mapping requests resulting from the **mmap** or **shmat** subroutines. The logical file system creates the virtual memory object (if it does not already exist) and increments the object's use count.

## Execution Environment

The `vnop_map` entry point can be called from the process environment only.

## Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure.

### Related reference:

“`vnop_fsync`, `vnop_fsync_range` Entry Points” on page 660

### Related information:

`shmat` subroutine

Virtual File System Kernel Extensions Overview

## `vnop_map_lloff` Entry Point

### Purpose

Announces intention to map a file.

### Syntax

```
int vnop_map_lloff (vp, addr, length, offset, mflags, fflags, crp) struct vnode *vp; caddr_t addr; offset_t length; offset_t offset; int mflags; int fflags; struct ucred *crp;
```

### Parameters

Item	Description
<i>vp</i>	Points to the vnode to be queried.
<i>addr</i>	Unused.
<i>length</i>	Specifies the length of the mapping request.
<i>offset</i>	Specifies the starting offset for the map request.
<i>mflags</i>	Specifies the mapping flags.
<i>fflags</i>	Specifies the file flags.
<i>crp</i>	Specifies user's credentials.

### Description

The `vnop_map_lloff` entry point is used to tell the file system that the file is going to be accessed by memory mapped loads and stores. The file system should fail the request if it does not support memory mapping. This interface allows applications to specify starting offsets that are larger than 2 gigabytes.

File systems that do not define `GFS_VERSION421` in their `gfs` flags do not need to supply a `vnop_map_lloff` entry point.

## Execution Environment

The `vnop_map_lloff` entry point can be called from the process environment only.

## Return Values



Item	Description
Zero	Indicates a successful operation.
Nonzero	Indicates that the operation failed; return values should be chosen from the <code>/usr/include/sys/errno.h</code> file.

#### Related information:

shmat subroutine  
mmap subroutine  
Virtual File System Kernel Extensions Overview

## vnop\_memcntl Entry Point Purpose

Changes or queries the physical attachment of a file.

### Syntax

```
#include <sys/vnode.h>
#include <sys/fcntl.h>
```

```
int vnop_memcntl (vnode, cmd, arg, crp)
struct vnode * vnode;
int cmd;
void * arg;
struct ucred * crp;
```

### Parameters

Item	Description
<i>vnode</i>	Points to the virtual node of the file
<i>cmd</i>	Specifies the operation to be performed. The <i>cmd</i> parameter can be one of the following values: <ul style="list-style-type: none"> <li>• F_ATTACH</li> <li>• F_DETACH</li> <li>• F_ATTINFO</li> </ul>
<i>arg</i>	Points to a structure containing the <code>attach_desc_t</code> , <code>detach_desc_t</code> or <code>attinfo_desc_t</code> information according to the specified <i>cmd</i> parameter. <pre>F_ATTACH attach_desc_t F_DETACH detach_desc_t F_ATTINFO attinfo_desc_t</pre>
<i>crp</i>	Points to the <code>cred</code> structure. This structure contains data that the file system can use to validate access permission.

### Description

The `vnop_memcntl` entry point requests memory attachment operations as specified by the *cmd* parameter. The *cmd* parameter determines the *arg* structure.

### Execution Environment

The `vnop_memcntl` entry point can be called from the process environment only.

### Return Values

Item	Description
0	Success.
non-zero	Failure.

#### Related information:

Workload management

## **vnode\_mkdir** Entry Point Purpose

Creates a directory.

### Syntax

```
int vnode_mkdir ( dp, name, mode, crp)
struct vnode *dp;
caddr_t name;
int mode;
struct ucred *crp;
```

### Parameters

Item	Description
<i>dp</i>	Points to the virtual node (v-node) of the parent directory of a new directory. This v-node is held for the duration of the entry point.
<i>name</i>	Specifies the name of a new directory.
<i>mode</i>	Specifies the permission modes of a new directory.
<i>crp</i>	Points to the <b>ucred</b> structure. This structure contains data that the file system can use to validate access permission.

### Description

The **vnode\_mkdir** entry point is invoked by the logical file system as the result of the **mkdir** subroutine. The **vnode\_mkdir** entry point is expected to create the named directory in the parent directory associated with the *dp* parameter. The logical file system ensures that the *dp* parameter does not reside on a read-only file system.

### Execution Environment

The **vnode\_mkdir** entry point can be called from the process environment only.

### Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the **/usr/include/sys/errno.h** file to indicate failure.

#### Related information:

mkdir subroutine

Virtual File System Overview

Logical File System Overview

## vnop\_mknod Entry Point

### Purpose

Creates a special file.

### Syntax

```
int vnop_mknod (dvp, name, mode, dev, crp)
struct vnode * dvp;
caddr_t * name;
int mode;
dev_t dev;
struct ucred * crp;
```

### Parameters

Item	Description
<i>dvp</i>	Points to the virtual node (v-node) for the directory to contain the new file. This v-node is held for the duration of the <b>vnop_mknod</b> entry point.
<i>name</i>	Specifies the name of a new file.
<i>mode</i>	Identifies the integer mode that indicates the type of file and its permissions.
<i>dev</i>	Identifies an integer device number.
<i>crp</i>	Points to the <b>cred</b> structure. This structure contains data that applications can use to validate access permission.

### Description

The **vnop\_mknod** entry point is invoked by the logical file system as the result of a **mknod** subroutine. The underlying file system is expected to create a new file in the given directory. The file type bits of the *mode* parameter indicate the type of file (regular, character special, or block special) to be created. If a special file is to be created, the *dev* parameter indicates the device number of the new special file.

The logical file system verifies that the *dvp* parameter does not reside in a read-only file system.

### Execution Environment

The **vnop\_mknod** entry point can be called from the process environment only.

### Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure.

#### Related information:

mknod subroutine

Virtual File System Overview

Logical File System Overview

## vnop\_open Entry Point

### Purpose

Requests that a file be opened for reading or writing.

## Syntax

```
int vnop_open (vp, flag, ext, vinfop, crp)
struct vnode * vp;
int flag;
caddr_t ext;
caddr_t vinfop;
struct ucred * crp;
```

## Parameters

Item	Description
<i>vp</i>	Points to the virtual node (v-node) associated with the desired file. The v-node is held for the duration of the open process.
<i>flag</i>	Specifies the type of access. Access modes are defined in the <code>/usr/include/sys/fcntl.h</code> file. <b>Note:</b> The <code>vnop_open</code> entry point does not use the FCREAT mode.
<i>ext</i>	Points to external data. This parameter is used if the subroutine is opening a device.
<i>vinfop</i>	This parameter is not currently used.
<i>crp</i>	Points to the <code>cred</code> structure. This structure contains data that the file system can use to validate access permission.

## Description

The `vnop_open` entry point is called to initiate a process access to a v-node and its underlying file system object. The operation of the `vnop_open` entry point varies between virtual file system (VFS) implementations. A successful `vnop_open` entry point must leave a v-node count of at least 1.

The logical file system ensures that the process is not requesting write access (with the FWRITE or FTRUNC mode) to a read-only file system.

## Execution Environment

The `vnop_open` entry point can be called from the process environment only.

## Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure.

### Related reference:

“`vnop_close` Entry Point” on page 655

### Related information:

open subroutine

Virtual File System Overview

## `vnop_rdwr`, `vnop_rdwr_attr` Entry Points

### Purpose

Requests file I/O.

### Syntax

```
int vnop_rdwr (vp, op, flags, uiop, ext, vinfo, vattrp, crp)
struct vnode * vp;
```

```

enum uio_rw op;
int flags;
struct uio * uiop;
int ext;
caddr_t vinfo;
struct vattr * vattrp;
struct ucred * crp;

int vnop_rdwr_attr (vp, op, flags, uiop, ext, vinfo, vpre, vpost, crp)
struct vnode * vp;
enum uio_rw op;
long flags;
struct uio * uiop;
ext_t ext;
caddr_t vinfo;
struct vattr * vpre;
struct vattr * vpost;
struct ucred * crp;

```

## Parameters

Item	Description
<i>vp</i>	Points to the virtual node (v-node) of the file.
<i>op</i>	Specifies a number that indicates a read or write operation. This parameter has a value of either <b>UIO_READ</b> or <b>UIO_WRITE</b> . These values are found in the <code>/usr/include/sys/uio.h</code> file.
<i>flags</i>	Identifies flags from the open file structure.
<i>uiop</i>	Points to a <b>uio</b> structure. This structure describes the count, data buffer, and other I/O information.
<i>ext</i>	Provides an extension for special purposes. Its use and meaning are specific to virtual file systems, and it is usually ignored except for devices.
<i>vinfo</i>	This parameter is currently not used.
<i>vattrp</i>	Points to a <b>vattr</b> structure. If this pointer is <b>NULL</b> , no action is required of the file system implementation. If it is not <b>NULL</b> , the attributes of the file specified by the <i>vp</i> parameter are returned at the address passed in the <i>vattrp</i> parameter.
<i>vpre</i>	Points to an attributes structure for pre-operation attributes.
<i>vpost</i>	Points to an attributes structure for post-operation attributes.
<i>crp</i>	Points to the <b>cred</b> structure. This structure contains data that the file system can use to validate access permission.

## Description

The **vnop\_rdwr** entry point is used to request that data to be read or written from an object represented by a v-node. The **vnop\_rdwr** entry point does the indicated data transfer and sets the number of bytes *not* transferred in the `uio_resid` field. This field is 0 (zero) on successful completion.

The **vnop\_rdwr\_attr** kernel service performs the same function as the **vnop\_rdwr** kernel service but also allows the caller to retrieve attributes of the object either before the I/O, after or both.

## Execution Environment

The **vnop\_rdwr** and **vnop\_rdwr\_attr** entry points can be called from the process environment only.

## Return Values

Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure. The **vnop\_rdwr** entry point returns an error code if an operation did not transfer all the data requested. The only exception is if an end of file is reached on a read request. In this case, the operation still returns 0.

**Related reference:**

“vnode\_create Entry Point” on page 656

“vnode\_open Entry Point” on page 675

#### Related information:

Virtual File System Kernel Extensions Overview

## vnode\_readdir Entry Point

### Purpose

Reads directory entries in standard format.

### Syntax

```
int vnode_readdir ( vp, uiop, crp)
struct vnode *vp;
struct uiop *uiop;
struct ucred *crp;
```

### Parameters

Item	Description
<i>vp</i>	Points to the virtual node (v-node) of the directory.
<i>uiop</i>	Points to the <b>uiop</b> structure that describes the data area into which to put the block of <b>dirent</b> structures. The starting directory offset is found in the <i>uiop-&gt;uio_offset</i> field and the size of the buffer area is found in the <i>uiop-&gt;uio_resid</i> field.
<i>crp</i>	Points to the <b>cred</b> structure. This structure contains data that the file system can use to validate access permission.

### Description

The **vnode\_readdir** entry point is used to access directory entries in a standard way. These directories should be returned as an array of **dirent** structures. The `/usr/include/sys/dir.h` file contains the definition of a **dirent** structure.

The **vnode\_readdir** entry point does the following:

- Copies a block of directory entries into the buffer specified by the *uiop* parameter.
- Sets the *uiop->uio\_resid* field to indicate the number of bytes read.

The End-of-file character should be indicated by not reading any bytes (not by a partial read). This provides directories with the ability to have some hidden information in each block.

The virtual file system-specific implementation is also responsible for setting the *uio\_offset* field to the offset of the next whole block to be read.

#### Notes:

- If the call is meant for a JFS2 filesystem, extra processing is needed to avoid duplicate entries being returned in the user data area. The caller can check the VFS type from the directory vnode.
- The caller must allocate a two-element array of type **struct iovec** to pass with the **uiop** structure. The first element is initialized to point to the user data area to receive the **dirent** structures. If the file pointer of the directory has a non-NULL *f\_vinfo* field, the second **iovec** element is initialized to point to the *f\_vinfo* field and the length is set to 0; the number of elements in the **uiop** structure is set to 2. If the *f\_vinfo* field is NULL, then the number of elements in the **uiop** structure is set to 1 and the second **iovec** element remains uninitialized.
- If the caller does not have access to the directory file pointer, a **dirent** structure can be allocated in place of the *f\_vinfo* field. The caller must not change this allocated structure between calls to the **vnode\_readdir** entry point.

## Execution Environment

The `vnop_readdir` entry point can be called from the process environment only.

## Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure.

### Related information:

`readdir` subroutine

Virtual File System Overview

Logical File System Overview

## `vnop_readdir_eofp` Entry Point

### Purpose

Returns directory entries.

### Syntax

```
int vnop_readdir_eofp (vp, uiop, eofp, crp) struct vnode *vp; struct uiop *uiop; int *eofp; struct ucred *crp;
```

### Parameters

Item	Description
<i>vp</i>	Points to the directory vnode to be processed.
<i>uiop</i>	Points to the uiop structure describing the user's buffer.
<i>eofp</i>	Points to a word that places the eop structure.
<i>crp</i>	Specifies user's credentials.

### Description

The `vnop_readdir_eofp` entry point is used to read directory entries. It is similar to `vnop_readdir` except that it takes the additional parameter, `eofp`. The location pointed to by the `eofp` parameter should be set to 1 if the `readdir` request reached the end of the directory. Otherwise, it should be set to 0.

File systems that do not define `GFS_VERSION421` in their `gfs` flags do not need to supply a `vnop_readdir_eofp` entry point.

**Note:** If the call is meant for a JFS2 file system, extra processing is needed to avoid duplicate entries being returned in the user data area, similar to the `vnop_readdir` entry point.

## Execution Environment

The `vnop_readdir_eofp` entry point can be called from the process environment only.

## Return Values

Item	Description
Zero	Indicates a successful operation.
Nonzero	Indicates that the operation failed; return values should be chosen from the <code>/usr/include/sys/errno.h</code> file.

#### Related reference:

“`vnop_readdir` Entry Point” on page 678

#### Related information:

`readdir` subroutine

Virtual File System Overview

Logical File System Overview

## `vnop_readlink` Entry Point

### Purpose

Reads the contents of a symbolic link.

### Syntax

```
int vnop_readlink ( vp, uio, crp)
struct vnode *vp;
struct uio *uio;
struct ucred *crp;
```

### Parameters

Item	Description
<i>vp</i>	Points to a virtual node (v-node) structure. The <code>vnop_readlink</code> entry point holds this v-node for the duration of the routine.
<i>uio</i>	Points to a <code>uio</code> structure. This structure contains the information required to read the link. In addition, it contains the return buffer for the <code>vnop_readlink</code> entry point.
<i>crp</i>	Points to the <code>cred</code> structure. This structure contains data that the file system can use to validate access permission.

### Description

The `vnop_readlink` entry point is used by the logical file system to get the contents of a symbolic link, if the file system supports symbolic links. The logical file system finds the v-node (virtual node) for the symbolic link, so this routine simply reads the data blocks for the symbol link.

### Execution Environment

The `vnop_readlink` entry point can be called from the process environment only.

### Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure.

#### Related reference:

“`uio` Structure” on page 639

#### Related information:

Virtual File System Kernel Extensions Overview

Logical File System Overview



## **vnode\_rele Entry Point**

### **Purpose**

Releases a reference to a virtual node (v-node).

### **Syntax**

```
int vnode_rele ( vp, )
struct vnode *vp;
```

### **Parameter**

Item	Description
<i>vp</i>	Points to the v-node.

### **Description**

The **vnode\_rele** entry point is used by the logical file system to release the object associated with a v-node. If the object was the last reference to the v-node, the **vnode\_rele** entry point then calls the **vn\_free** kernel service to deallocate the v-node.

If the virtual file system (VFS) was unmounted while there were open files, the logical file system sets the **VFS\_UNMOUNTING** flag in the **vfs** structure. If the flag is set and the v-node to be released is the last v-node on the chain of the **vfs** structure, then the virtual file system must be deallocated with the **vnode\_rele** entry point.

### **Execution Environment**

The **vnode\_rele** entry point can be called from the process environment only.

### **Return Values**

Item	Description
0	Indicates success.

Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure.

#### **Related reference:**

“vn\_free Kernel Service” on page 580

#### **Related information:**

Virtual File System Overview

Logical File System Overview

## **vnode\_remove Entry Point**

### **Purpose**

Unlinks a file or directory.

### **Syntax**

```
int vnode_remove ( vp, dvp, name, crp)
struct vnode *vp;
```

```

struct vnode *dvp;
char *name;
struct ucred *crp;

```

## Parameters

Item	Description
<i>vp</i>	Points to a virtual node (v-node). The v-node indicates which file to remove and is held over the duration of the <b>vnode_remove</b> entry point.
<i>dvp</i>	Points to the v-node of the parent directory. This directory contains the file to be removed. The directory's v-node is held for the duration of the <b>vnode_remove</b> entry point.
<i>name</i>	Identifies the name of the file.
<i>crp</i>	Points to the <b>cred</b> structure. This structure contains data that the file system can use to validate access permission.

## Description

The **vnode\_remove** entry point is called by the logical file system to remove a directory entry (or link) as the result of a call to the **unlink** subroutine.

The logical file system assumes that the **vnode\_remove** entry point calls the **vnode\_rele** entry point. If the link is the last reference to the file in the file system, the disk resources that the file is using are released.

The logical file system ensures that the directory specified by the *dvp* parameter does not reside in a read-only file system.

## Execution Environment

The **vnode\_remove** entry point can be called from the process environment only.

## Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure.

### Related reference:

“vnode\_rele Entry Point” on page 681

### Related information:

unlink subroutine

Virtual File System Overview

## vnode\_rename Entry Point

### Purpose

Renames a file or directory.

### Syntax

```

int vnode_rename (srcvp, srcdvp, oldname, destvp, destdvp, newname, crp)
struct vnode * srcvp;
struct vnode * srcdvp;
char * oldname;
struct vnode * destvp;

```

```

struct vnode * destdvp;
char * newname;
struct ucred * crp;

```

## Parameters

Item	Description
<i>srcvp</i>	Points to the virtual node (v-node) of the object to rename.
<i>srcdvp</i>	Points to the v-node of the directory where the <i>srcvp</i> parameter resides. The parent directory for the old and new object can be the same.
<i>oldname</i>	Identifies the old name of the object.
<i>destdvp</i>	Points to the v-node of the new object. This pointer is used only if the new object exists. Otherwise, this parameter is the null character.
<i>destdop</i>	Points to the parent directory of the new object. The parent directory for the new and old objects can be the same.
<i>newname</i>	Points to the new name of the object.
<i>crp</i>	Points to the <b>cred</b> structure. This structure contains data that applications can use to validate access permission.

## Description

The **vnop\_rename** entry point is invoked by the logical file system to rename a file or directory. This entry point provides the following renaming actions:

- Renames an old object to a new object that exists in a different parent directory.
- Renames an old object to a new object that does not exist in a different parent directory.
- Renames an old object to a new object that exists in the same parent directory.
- Renames an old object to a new object that does not exist in the same parent directory.

To ensure that this entry point routine executes correctly, the logical file system guarantees the following:

- File names are not renamed across file systems.
- The old and new objects (if specified) are not the same.
- The old and new parent directories are of the same type of v-node.

## Execution Environment

The **vnop\_rename** entry point can be called from the process environment only.

## Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure.

### Related information:

rename subroutine

Virtual File System Overview

Logical File System Overview

## vnop\_revoke Entry Point

### Purpose

Revokes all access to an object.

## Syntax

```
int vnop_revoke (vp, cmd, flag, vinfop, crp)
struct vnode * vp;
int cmd;
int flag;
caddr_t vinfop;
struct ucred * crp;
```

## Parameters

Item	Description
<i>vp</i>	Points to the virtual node (v-node) containing the object.
<i>cmd</i>	Indicates whether the calling process holds the file open. This parameter takes the following values:  0           The process did not have the file open.  1           The process had the file open.  2           The process had the file open and the reference count in the file structure was greater than 1.
<i>flag</i>	Identifies the flags from the <b>file</b> structure.
<i>vinfop</i>	This parameter is currently unused.
<i>crp</i>	Points to the <b>cred</b> structure. This structure contains data that the file system can use to validate access permission.

## Description

The **vnop\_revoke** entry point is called to revoke further access to an object.

## Execution Environment

The **vnop\_revoke** entry point can be called from the process environment only.

## Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the **/usr/include/sys/errno.h** file to indicate failure.

### Related information:

frevoke subroutine

revoke subroutine

Virtual File System Overview

## vnop\_rmdir Entry Point

### Purpose

Removes a directory.

## Syntax

```
int vnop_rmdir ( vp, dp, pname, crp)
struct vnode *vp;
struct vnode *dp;
char *pname;
struct ucred *crp;
```

## Parameters

Item	Description
<i>vp</i>	Points to the virtual node (v-node) of the directory.
<i>dp</i>	Points to the parent of the directory to remove.
<i>pname</i>	Points to the name of the directory to remove.
<i>crp</i>	Points to the <b>cred</b> structure. This structure contains data that the file system can use to validate access permission.

## Description

The **vnop\_rmdir** entry point is invoked by the logical file system to remove a directory object. To remove a directory, the directory must be empty (except for the current and parent directories). Before removing the directory, the logical file system ensures the following:

- The *vp* parameter is a directory.
- The *vp* parameter is not the root of a virtual file system.
- The *vp* parameter is not the current directory.
- The *dp* parameter does not reside on a read-only file system.

**Note:** The *vp* and *dp* parameters' v-nodes (virtual nodes) are held for the duration of the routine.

## Execution Environment

The **vnop\_rmdir** entry point can be called from the process environment only.

## Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure.

### Related information:

[rmdir subroutine](#)

[Virtual File System Overview](#)

[Logical File System Overview](#)

## vnop\_seek Entry Point

### Purpose

Validates file offsets.

### Syntax

```
int vnop_seek (vp, offsetp, crp)
struct vnode * vp;
offset_t * offp;
struct ucred * crp;
```

## Parameters

Item	Description
<i>vp</i>	Points to the virtual node (vnode) of the file.
<i>offp</i>	Points to the location of the new offset to validate.
<i>crp</i>	Points to the user's credential.

## Description

The **vnode\_seek** entry point is called by the logical file system to validate a new offset that is computed by the **lseek**, **llseek**, and **lseek64** subroutines. The file system implementation must check the offset that is pointed to by the *offp* parameter and, if it is acceptable for the file, return zero. If the offset is not acceptable, the routine must return a non-zero value. **EINVAL** is the suggested error value for invalid offsets.

File system that do not want to do offset validation can simply return 0. File system that do not provide the **vnode\_seek** entry point has a maximum offset of **OFF\_MAX** (2 gigabytes minus 1) enforced by the logical file system.

## Execution Environment

The **vnode\_seek** entry point is to be called from the process environment only.

## Return Values

Item	Description
0	Indicates success.
Non-zero	Return values are returned the <code>/usr/include/sys/errno.h</code> file to indicate failure.

### Related information:

**lseek**, **llseek**, and **lseek64**

## vnode\_select Entry Point Purpose

Polls a virtual node (v-node) for immediate I/O.

## Syntax

```
int vnode_select (vp, correl, e, re, notify, vinfo, crp)
struct vnode * vp;
int correl;
int e;
int re;
int (* notify)();
caddr_t vinfo;
struct ucred * crp;
```

## Parameters

Item	Description
<i>vp</i>	Points to the v-node to be polled.
<i>correl</i>	Specifies the ID used for correlation in the <b>selnotify</b> kernel service.
<i>e</i>	Identifies the requested event.
<i>re</i>	Returns an events list. If the v-node is ready for immediate I/O, this field should be set to indicate the requested event is ready.
<i>notify</i>	Specifies the subroutine to call when the event occurs. This parameter is for nested polls.
<i>vinfo</i>	Is currently unused.
<i>crp</i>	Points to the <b>cred</b> structure. This structure contains data that the file system can use to validate access permission.

## Description

The **vnop\_select** entry point is invoked by the logical file system to poll a v-node to determine if it is immediately ready for I/O. This entry point is used to implement the **select** and **poll** subroutines.

File system implementation can support constructs, such as devices or pipes, that support the select semantics. The **fp\_select** kernel service provides more information about select and poll requests.

## Execution Environment

The **vnop\_select** entry point can be called from the process environment only.

## Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure.

### Related reference:

“**fp\_select** Kernel Service” on page 167

### Related information:

select subroutine

Virtual File System Kernel Extensions Overview

## vnop\_setacl Entry Point

### Purpose

Sets the access control list (ACL) for a file.

### Syntax

```
#include <sys/acl.h>
```

```
int vnop_setacl ( vp, uiop, crp)
struct vnode *vp;
struct uio *uiop;
struct ucred *crp;
```

### Description

The **vnop\_setacl** entry point is used by the logical file system to set the access control list (ACL) on a file.

## Parameters

Item	Description
<i>vp</i>	Specifies the virtual node (v-node) of the file system object.
<i>uiop</i>	Specifies the <b>uio</b> structure that defines the storage for the call arguments.
<i>crp</i>	Points to the <b>cred</b> structure. This structure contains data that the file system can use to validate access permission.

## Execution Environment

The **vnop\_setacl** entry point can be called from the process environment only.

## Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the **/usr/include/sys/errno.h** file to indicate failure. Valid values include:

Item	Description
ENOSPC	Indicates that the space cannot be allocated to hold the new ACL information.
EPERM	Indicates that the effective user ID of the process is not the owner of the file and the process is not privileged.

### Related information:

chacl subroutine

statacl subroutine

Virtual File System Overview

## vnop\_setattr Entry Point

### Purpose

Sets attributes of a file.

### Syntax

```
int vnop_setattr (vp, cmd, arg1, arg2, arg3, crp)
struct vnode * vp;
int cmd;
int arg1;
int arg2;
int arg3;
struct ucred * crp;
```

### Description

The **vnop\_setattr** entry point is used by the logical file system to set the attributes of a file. This entry point is used to implement the **chmod**, **chownx**, and **utime** subroutines.

The values that the *arg* parameters take depend on the value of the *cmd* parameter. The **vnop\_setattr** entry point accepts the following *cmd* values and *arg* parameters:



## Possible cmd Values for the vnop\_setattr Entry Point

Command	V_OWN	V_UTIME	V_MODE
<i>arg1</i>	<b>int</b> <i>flag</i> ;	<b>int</b> <i>flag</i> ;	<b>int</b> <i>mode</i> ;
<i>arg2</i>	<b>int</b> <i>uid</i> ;	<b>timestruc_t</b> * <i>atime</i> ;	Unused
<i>arg3</i>	<b>int</b> <i>gid</i> ;	<b>timestruc_t</b> * <i>mtime</i> ;	Unused

**Note:** For **V\_UTIME**, if *arg2* or *arg3* is **NULL**, then the corresponding time field, *atime* and *mtime*, of the file should be left unchanged.

## Parameters

Item	Description
<i>vp</i>	Points to the virtual node (v-node) of the file.
<i>cmd</i>	Defines the setting operation. This parameter takes the following values: <ul style="list-style-type: none"> <li><b>V_OWN</b> Sets the user ID (UID) and group ID (GID) to the UID and GID values of the new file owner. The <i>flag</i> argument indicates which ID is affected.</li> <li><b>V_UTIME</b> Sets the access and modification time for the new file. If the <i>flag</i> parameter has the value of <b>T_SETTIME</b>, then the specific values have not been provided and the access and modification times of the object should be set to current system time. If the <b>T_SETTIME</b> value is not specified, the values are specified by the <i>atime</i> and <i>mtime</i> variables.</li> <li><b>V_MODE</b> Sets the file mode.  The <code>/usr/include/sys/vattr.h</code> file contains the definitions for the three command values.</li> </ul>
<i>arg1, arg2, arg3</i>	Specify the command arguments. The values of the command arguments depend on which command calls the <b>vnop_setattr</b> entry point.
<i>crp</i>	Points to the <b>cred</b> structure. This structure contains data that the file system can use to validate access permission.

## Execution Environment

The **vnop\_setattr** entry point can be called from the process environment only.

## Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure.

### Related information:

chmod subroutine

Understanding Virtual Nodes (V-nodes)

## vnop\_setxacl Entry Point

### Purpose

Sets the access control list (ACL) for a file system object. This is an advanced interface compared to **vnop\_setacl** and provides for ACL-type-based operations.

## Syntax

```
#include <sys/acl.h>
int vnop_setxacl (vp, ctl_flags, acl_type, uiop, mode_info, crp)

struct vnode    *vp;
uint64_t       ctl_flags;
acl_type_t     acl_type;
struct uio     *uiop;
mode_t         mode_info;
struct ucred   *crp;
```

## Description

The `vnop_setxacl` entry point sets the access control list (ACL) on a file. It is an advanced version of `vnop_setacl` interface and provides for ACL-type-based operations. This interface can also be used to manage special bits in mode word (such as SUID, SGID and SVTX) in case the ACL type does not support these bits through ACL.

## Parameters

Item	Description
<code>vp</code>	Specifies the virtual node (v-node) of the file system object for which the ACL needs to be set.
<code>acl_type</code>	Specifies the ACL type of the ACL information that needs to be set for the file system object. <b>Note:</b> If the underlying physical file system does not support the ACL type being requested, the system could return an error.
<code>acl_len</code>	Pointer to a <i>length</i> variable. The space pointed to is used as an input, as well as output, parameter. As input, the value will indicate the size of buffer <i>uiop</i> . When the call returns, this space holds the actual length of the ACL (true for when the call is successful or when the call fails with <b>errno</b> set to <b>ENOSPC</b> ).
<code>ctl_flags</code>	This 64-bit bit mask provides for control over the ACL setting and for any future variations in the interface. The following flag values have been defined:  <b>SET_MODE_S_BITS</b> Indicates that the <i>mode_info</i> value is set by the caller and the ACL put operation must consider this value to complete the ACL put operation.  <b>SET_ACL</b> Indicates that the ACL arguments point to valid ACL data that must be considered while the ACL put operation is being performed. <b>Note:</b> Both of the preceding values can be specified by the caller by ORing the two masks.
<code>uiop</code>	Specifies the <b>uio</b> structure that defines the storage for the call arguments.
<code>mode_info</code>	This value indicates any mode word information that needs to be set for the file system object as part of this ACL put operation. When mode bits are altered by specifying the <b>SET_MODE_S_BITS</b> flag (in <i>ctl_flags</i> ), the entire ACL put operation will fail if the caller does not have the required privileges.
<code>crp</code>	Points to the <b>cred</b> structure. This structure contains data that the file system can use to validate access permission.

## Execution Environment

The `vnop_setxacl` entry point can be called from the process environment only.

## Return Values

Upon successful completion, the `vnop_setxacl` entry point returns 0. Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure.

Item	Description
EPERM	Indicates that the effective user ID of the process is not authorized to change the ACL on the specified file system object.
EINVAL	Invalid operation. File system might not support the ACL type being set.

**Note:** This list of error numbers is not complete and is dependent on the particular physical file system implementation supporting the ACL.

**Related reference:**

“vnode\_setacl Entry Point” on page 687

**Related information:**

statacl subroutine

Logical File System Overview

## vnode\_strategy Entry Point Purpose

Accesses blocks of a file.

### Syntax

```
int vnode_strategy ( vp, bp, crp)
struct vnode *vp;
struct buf *bp;
struct ucred *crp;
```

### Description

**Note:** The `vnode_strategy` entry point is not implemented in Version 3.2 of the operating system.

The `vnode_strategy` entry point accesses blocks of a file. This entry point is intended to provide a block-oriented interface for servers for efficiency in paging.

### Parameters

Item	Description
<i>vp</i>	Points to the virtual node (v-node) of the file.
<i>bp</i>	Points to a <code>buf</code> structure that describes the buffer.
<i>crp</i>	Points to the <code>cred</code> structure. This structure contains data that applications can use to validate access permission.

### Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure.

**Related reference:**

“buf Structure” on page 615

**Related information:**

Virtual File System Overview

Virtual File System Kernel Extensions Overview

## vnop\_symlink Entry Point

### Purpose

Creates a symbolic link.

### Syntax

```
int vnop_symlink ( vp, linkname, target, crp)
struct vnode *vp;
char *linkname;
char *target;
struct ucred *crp;
```

### Description

The `vnop_symlink` entry point is called by the logical file system to create a symbolic link. The path name specified by the `linkname` parameter is the name of the new symbolic link. This symbolic link points to the object named by the `target` parameter.

### Parameters

Item	Description
<i>vp</i>	Points to the virtual node (v-node) of the parent directory where the link is created.
<i>linkname</i>	Points to the name of the new symbolic link. The logical file system guarantees that the new link does not already exist.
<i>target</i>	Points to the name of the object to which the symbolic link points. This name need not be a fully qualified path name or even an existing object.
<i>crp</i>	Points to the <code>ucred</code> structure. This structure contains data that the file system can use to validate access permission.

### Execution Environment

The `vnop_symlink` entry point can be called from the process environment only.

### Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure.

#### Related information:

symlink subroutine

Virtual File System Overview

Logical File System Overview

## vnop\_unmap Entry Point

### Purpose

Unmaps a file.

## Syntax

```
int vnop_unmap ( vp, flag, crp)
struct vnode *vp;
ulong flag;
struct ucred *crp;
```

## Description

The **vnop\_unmap** entry point is called by the logical file system to unmap a file. When this entry point routine completes successfully, the use count for the memory object should be decremented and (if the use count went to 0) the memory object should be destroyed. The file system implementation is required to perform only those operations that are unique to the file system. The logical file system handles virtual-memory management operations.

## Parameters

Item	Description
<i>vp</i>	Points to the v-node (virtual node) of the file.
<i>flag</i>	Indicates how the file was mapped. This flag takes the following values: <b>SHM_RDONLY</b> The virtual memory object is read-only. <b>SHM_COPY</b> The virtual memory object is copy-on-write.
<i>crp</i>	Points to the <b>cred</b> structure. This structure contains data that the file system can use to validate access permission.

## Execution Environment

The **vnop\_unmap** entry point can be called from the process environment only.

## Return Values

Item	Description
0	Indicates success.

Nonzero return values are returned from the `/usr/include/sys/errno.h` file to indicate failure.

### Related information:

Virtual File System Overview

Virtual File System Kernel Extensions Overview

Logical File System Overview



---

## Notices

This information was developed for products and services offered in the US.

IBM may not offer the products, services, or features discussed in this document in other countries. Consult your local IBM representative for information on the products and services currently available in your area. Any reference to an IBM product, program, or service is not intended to state or imply that only that IBM product, program, or service may be used. Any functionally equivalent product, program, or service that does not infringe any IBM intellectual property right may be used instead. However, it is the user's responsibility to evaluate and verify the operation of any non-IBM product, program, or service.

IBM may have patents or pending patent applications covering subject matter described in this document. The furnishing of this document does not grant you any license to these patents. You can send license inquiries, in writing, to:

*IBM Director of Licensing  
IBM Corporation  
North Castle Drive, MD-NC119  
Armonk, NY 10504-1785  
US*

For license inquiries regarding double-byte character set (DBCS) information, contact the IBM Intellectual Property Department in your country or send inquiries, in writing, to:

*Intellectual Property Licensing  
Legal and Intellectual Property Law  
IBM Japan Ltd.  
19-21, Nihonbashi-Hakozakicho, Chuo-ku  
Tokyo 103-8510, Japan*

INTERNATIONAL BUSINESS MACHINES CORPORATION PROVIDES THIS PUBLICATION "AS IS" WITHOUT WARRANTY OF ANY KIND, EITHER EXPRESS OR IMPLIED, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF NON-INFRINGEMENT, MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE. Some jurisdictions do not allow disclaimer of express or implied warranties in certain transactions, therefore, this statement may not apply to you.

This information could include technical inaccuracies or typographical errors. Changes are periodically made to the information herein; these changes will be incorporated in new editions of the publication. IBM may make improvements and/or changes in the product(s) and/or the program(s) described in this publication at any time without notice.

Any references in this information to non-IBM websites are provided for convenience only and do not in any manner serve as an endorsement of those websites. The materials at those websites are not part of the materials for this IBM product and use of those websites is at your own risk.

IBM may use or distribute any of the information you provide in any way it believes appropriate without incurring any obligation to you.

Licenseses of this program who wish to have information about it for the purpose of enabling: (i) the exchange of information between independently created programs and other programs (including this one) and (ii) the mutual use of the information which has been exchanged, should contact:

*IBM Director of Licensing  
IBM Corporation  
North Castle Drive, MD-NC119  
Armonk, NY 10504-1785  
US*

Such information may be available, subject to appropriate terms and conditions, including in some cases, payment of a fee.

The licensed program described in this document and all licensed material available for it are provided by IBM under terms of the IBM Customer Agreement, IBM International Program License Agreement or any equivalent agreement between us.

The performance data and client examples cited are presented for illustrative purposes only. Actual performance results may vary depending on specific configurations and operating conditions.

Information concerning non-IBM products was obtained from the suppliers of those products, their published announcements or other publicly available sources. IBM has not tested those products and cannot confirm the accuracy of performance, compatibility or any other claims related to non-IBM products. Questions on the capabilities of non-IBM products should be addressed to the suppliers of those products.

Statements regarding IBM's future direction or intent are subject to change or withdrawal without notice, and represent goals and objectives only.

All IBM prices shown are IBM's suggested retail prices, are current and are subject to change without notice. Dealer prices may vary.

This information is for planning purposes only. The information herein is subject to change before the products described become available.

This information contains examples of data and reports used in daily business operations. To illustrate them as completely as possible, the examples include the names of individuals, companies, brands, and products. All of these names are fictitious and any similarity to actual people or business enterprises is entirely coincidental.

#### COPYRIGHT LICENSE:

This information contains sample application programs in source language, which illustrate programming techniques on various operating platforms. You may copy, modify, and distribute these sample programs in any form without payment to IBM, for the purposes of developing, using, marketing or distributing application programs conforming to the application programming interface for the operating platform for which the sample programs are written. These examples have not been thoroughly tested under all conditions. IBM, therefore, cannot guarantee or imply reliability, serviceability, or function of these programs. The sample programs are provided "AS IS", without warranty of any kind. IBM shall not be liable for any damages arising out of your use of the sample programs.

Each copy or any portion of these sample programs or any derivative work must include a copyright notice as follows:

© (your company name) (year).

Portions of this code are derived from IBM Corp. Sample Programs.

© Copyright IBM Corp. \_enter the year or years\_.



---

## Privacy policy considerations

IBM Software products, including software as a service solutions, (“Software Offerings”) may use cookies or other technologies to collect product usage information, to help improve the end user experience, to tailor interactions with the end user or for other purposes. In many cases no personally identifiable information is collected by the Software Offerings. Some of our Software Offerings can help enable you to collect personally identifiable information. If this Software Offering uses cookies to collect personally identifiable information, specific information about this offering’s use of cookies is set forth below.

This Software Offering does not use cookies or other technologies to collect personally identifiable information.

If the configurations deployed for this Software Offering provide you as the customer the ability to collect personally identifiable information from end users via cookies and other technologies, you should seek your own legal advice about any laws applicable to such data collection, including any requirements for notice and consent.

For more information about the use of various technologies, including cookies, for these purposes, see IBM’s Privacy Policy at <http://www.ibm.com/privacy> and IBM’s Online Privacy Statement at <http://www.ibm.com/privacy/details> the section entitled “Cookies, Web Beacons and Other Technologies” and the “IBM Software Products and Software-as-a-Service Privacy Statement” at <http://www.ibm.com/software/info/product-privacy>.

---

## Trademarks

IBM, the IBM logo, and [ibm.com](http://www.ibm.com) are trademarks or registered trademarks of International Business Machines Corp., registered in many jurisdictions worldwide. Other product and service names might be trademarks of IBM or other companies. A current list of IBM trademarks is available on the web at [www.ibm.com/legal/copytrade.shtml](http://www.ibm.com/legal/copytrade.shtml).

UNIX is a registered trademark of The Open Group in the United States and other countries.



---

# Index

## Special characters

\_\_pag\_getid system call 399  
\_\_pag\_getname System Call 400  
\_\_pag\_getvalue system call 400  
\_\_pag\_setname System Call 401  
\_\_pag\_setvalue system call 402

## A

access control lists  
  retrieving 662, 664  
  setting 642, 687, 689  
acct\_add\_LL Kernel Service 2  
acct\_get\_projid Kernel Service 2  
acct\_get\_usage Kernel Service 3  
acct\_interval\_register Kernel Service 4  
acct\_interval\_unregister Kernel Service 4  
acct\_put Kernel Service 6  
acct\_zero\_LL Kernel Service 2  
add\_domain\_af kernel service 8  
add\_input\_type kernel service 8  
add\_netisr kernel service 10  
add\_netopt macro 10  
address families  
  adding 8  
  deleting 64  
  searching for 405  
address ranges  
  pinning 354, 408, 605  
  setting storage protect key for 561  
  unpinning 355, 520, 606  
address space  
  kernel memory  
    allocating 11  
    deallocating 12  
    mapping 11, 21  
    obtaining handles 13, 14, 15  
    releasing 20  
    unmapping 12  
advanced accounting  
  acct\_add\_LL Kernel Service 2  
  acct\_get\_projid Kernel Service 2  
  acct\_get\_usage Kernel Service 3  
  acct\_interval\_register Kernel Service 4  
  acct\_interval\_unregister Kernel Service 4  
  acct\_put Kernel Service 6  
  acct\_zero\_LL Kernel Service 2  
allocate memory  
  rmalloc 449  
allocated memory  
  freeing 611  
allocating memory  
  rmfree 450  
as\_att64 kernel service  
  described 11  
as\_det64 kernel service 12  
as\_geth kernel service 13  
as\_geth64 kernel service 14  
as\_getsrval64 kernel service 15  
as\_lw\_att64 Kernel Service 16

as\_lw\_det64 Kernel Service 18  
as\_lw\_pool\_init Kernel Service 19  
as\_puth64 kernel service 20  
as\_seth64 kernel service 21  
asynchronous processing  
  notify routine and 168  
asynchronous requests  
  registering 465  
attach-device queue management routine 22  
audit records  
  appending to 23  
  completing 24  
  initiating 24  
  writing 24  
audit\_svcbcopy kernel service 23  
audit\_svcfinis kernel service 24  
audit\_svcstart kernel service 24

## B

bawrite kernel service 26  
bdwrite kernel service 26  
bflush kernel service 27  
binding a process to a processor 28  
bindprocessor kernel service 28  
binval kernel service 29  
blkflush kernel service 30  
block I/O  
  buf headers  
    completion of 526  
    preparing 525  
  buf structures 615  
  calling 526  
  character I/O for blocks  
    performing 524  
  completion  
    waiting for 233  
  requests  
    completing 227  
block I/O buffer cache  
  assigning blocks 30  
  assigning buffer 181  
  buf structures 615  
  buffers  
    header address 186  
    purging block from 422  
  clearing 42  
  flushing 30  
  freeing 32  
  nonreclaimable blocks 29  
  read-ahead block 31  
  reading blocks into 30, 31  
  releasing 26  
  write-behind blocks 27  
  writing 36  
  writing contents asynchronously 26  
  zeroing-out 42  
blocked processes  
  clearing 510  
blocks  
  purging from buffer 422

- bread kernel service 30
- breada kernel service 31
- brlse kernel service 32
- bsr\_alloc Kernel Service 33
- bsr\_free Kernel Service 34
- bsr\_query Kernel Service 35
- buf headers
  - completion of 526
  - preparing 525
  - sending to a routine 527
- buf structures 615
- buffer cache 26
- buffers 185
  - allocating 186
  - determining status 187
  - freeing 426
  - freeing buffer lists 426
  - header address of 186
- bufx structure 617
- bus interrupt levels
  - disabling 219
  - enabling 239
  - resetting 237
- bwrite kernel service 36
- bytes
  - storing 478

## C

- caller's buffer
  - md\_restart\_block\_read 373
- callout table entries
  - registering changes in 492
- cancel pending timer requests 523
- cancel-queue-element queue management routine 37
- cascade processing 168
- cfgnadd kernel service 38
- cfgncb control block
  - adding 38
  - removing 40
- cfgncb kernel service 39
- cfgndel kernel service 40
- chan parameter 614
- channel numbers
  - finding 148
- character data
  - reading from device 631
- character device driver
  - character lists 619
  - clist structure 619
- character I/O
  - freeing buffers 185
  - getting buffer addresses 183
  - performing for blocks 524
  - placing character buffers 424
  - placing characters 425, 427
  - placing characters in list 423
  - retrieving a character 182
  - retrieving from buffers 531
  - retrieving last character 185
  - retrieving multiple characters 184
  - uio structures 639
  - writing to buffers 530
- character lists
  - removing first buffer 183
  - structure of 619
  - using 619

- check-parameters queue management routine 41
- close subroutine
  - device driver 620
- clrbuf kernel service 42
- clrjmx kernel service 43
- common\_reclock kernel service 43
- communication I/O device handler
  - opening 378
- communications device handlers
  - closing 379
  - transmitting data to 384
- compare\_and\_swap kernel service 45
- compare\_and\_swaplp kernel service 45
- configuration notification control block 39
- contexts
  - saving 468
- conventional locks
  - locking 343
- coprocessor\_user\_register kernel service 46
- coprocessor\_user\_unregister kernel service 47
- copyin kernel service 48
- copying to NVAM header
  - md\_restart\_block\_upd Kernel Service 374
- copyinstr kernel service 49
- copyout kernel service 50
- cpu\_speculation\_barrier Kernel service 51
- creatp kernel service 54
- cross-memory move
  - performing 609
- ctlinput function
  - invoking 404
- curtime kernel service 60

## D

- d\_align kernel service 61
- d\_alloc\_dmamem kernel service 61
- d\_cflush kernel service 62
- d\_free\_dmamem kernel service 76
- d\_map\_attr kernel service 79
- d\_map\_clear kernel service 80
- d\_map\_disable kernel service 80
- d\_map\_enable 81
- d\_map\_init kernel service 82
- d\_map\_init\_ext kernel service 83
- d\_map\_list kernel service 84
- d\_map\_page kernel service 86
- d\_map\_query kernel service 88
- d\_map\_slave 90
- d\_roundup kernel service 104
- d\_sync\_mem kernel service 105
- d\_unmap\_list kernel service 107
- d\_unmap\_page kernel service 108
- d\_unmap\_slave 107
- data
  - memory
    - moving to kernel global memory 608
  - moving
    - from kernel global memory 609
    - moving between VMO and buffer 551
    - retrieving a byte 179
    - sending to DLC 173
  - word
    - retrieving 180
- data blocks
  - moving 517
- ddclose entry point 620

- ddconfig entry point 621
- dddump entry point
  - calling 69
  - writing to a device 624
- ddioctl entry point 626
- ddmpx entry point 627
- ddopen entry point 629
- ddread entry point
  - reading data from a character device 631
- ddrevoke entry point 632
- ddselect entry point
  - occurring on a device 633
- ddselect routine
  - calling fp\_select kernel service 168
- ddstrategy entry point
  - block-oriented I/O 635
  - calling 70
- ddwrite entry point
  - writing to a character device 637
- de-allocate resource
  - d\_unmap\_slave 107
- deallocates resources
  - d\_map\_clear 80
  - d\_unmap\_list 107
- del\_domain\_af kernel service 64
- del\_input\_type kernel service 65
- del\_netisr kernel service 66
- delay kernel service 64
- destination addresses
  - locating 212
- devdump kernel service 69
- device driver 614
  - access
    - revoking 632
  - buf structures 615
  - character data
    - reading 631
  - closing 620
  - configuration data
    - requesting 621
  - configuring 621
  - data
    - writing 637
  - events
    - checking for 633
  - iodone kernel service 227
  - memory buffers 639
  - multiplexed
    - allocating channels 627
    - deallocating channels 627
  - performing block-oriented I/O 635
  - performing special operations 626
  - preparing for control functions 629
  - preparing for reading 629
  - preparing for writing 629
  - read logic
    - reads and writes 639
  - select logic
    - reads and writes 639
  - terminating 621
  - uio structures 639
- device driver entry points
  - ddclose 620
  - ddconfig
    - writing to a device 621
  - dddump
    - writing to a device 624
- device driver entry points *(continued)*
  - ddioctl 626
  - ddmpx 627
  - ddopen 629
  - ddread 631
  - ddrevoke 632
  - ddselect 633
  - ddstrategy 635
  - ddwrite 637
  - standard parameters 614
- device driver management
  - dddump entry point
    - calling 69
  - ddstrategy entry point
    - calling 70
  - device entry
    - status 75
  - disk driver tasks 230
  - dkstat structure 229
  - entry points
    - adding 71
    - deleting 73
    - function pointers 279
  - exception handlers
    - deleting system-wide 510
    - system-wide 506
  - exception information
    - retrieving 188
  - kernel object files
    - loading 280
    - unloading 283
  - notification routines
    - adding 417
    - deleting 419
  - poll request
    - support for 463
  - processes
    - blocking 509
    - clearing blocked 510
  - programmed I/O
    - exceptions caused by 411
  - registering asynchronous requests 465
  - registering notification routine 38
  - removing control blocks 40
  - select request
    - support for 463
  - statistics structures
    - registering 229
    - removal 232
  - symbol binding support 282
  - ttystat structure 229
  - u\_error fields 191
  - ut\_error field
    - setting 469
- device handlers
  - ending a start 383
  - pio\_assist kernel service 412
  - starting network ID on 382
- device numbers
  - finding 148
- device queue management
  - attachq kernel service support 22
  - control block structure 39
  - detchq kernel service support 68
  - queue elements
    - placing into queue 137
    - waiting for 583

device queue management (*continued*)

virtual interrupt handlers

defining 534

removing 533

device switch table

altering a 72

devices

select request on 167

devno parameter 614

devstrat kernel service 70

devswadd kernel service 71

devswchg kernel service 72

devswdel kernel service 73

devswqry kernel service 75

direct memory access 61

directories

creating 674

entries

reading 678

removing 684

renaming 682

unlinking 681

disable DMA

d\_map\_disable 80

disable\_lock kernel service 77

disablement\_checking\_resume Kernel Service 78

disablement\_checking\_suspend Kernel Service 78

disk driver support 230

dkstat structure 229

DLC kernel services

fp\_ioctl 152

fp\_open 157

fp\_write 173

trcgenkt 495

DLC management

channel

disabling 145

device manager

opening 157

file pointers

sending kernel data to 173

trace channels

recording events 495

transferring commands to 152

DMA

disable

d\_map\_disable 80

enable

d\_map\_enable 81

DMA management

address ranges

pinning 408, 605

unpinning 606

buffer cache

maintaining 104

cache

flushing 62

cache-line size 61

processor cache

flushing 537

DMA master devices

deallocates resources

d\_unmap\_page 108

mapping

d\_map\_page 86

DMA operations

allocates and initializes resources

d\_map\_init 82

dmp\_add kernel service 92

dmp\_comptext kernel service 93

dmp\_compspec kernel service 93

dmp\_context kernel service 100

dmp\_ct kernel service 100

dmp\_ctl kernel service 96

dmp\_del kernel service 100

dmp\_eaddr kernel service 100

dmp\_errbuf kernel service 100

dmp\_kernext kernel service 104

dmp\_mtrc kernel service 100

dmp\_pid kernel service 100

dmp\_systrace kernel service 100

dmp\_tid kernel service 100

dr\_reconfig system call 109

DTOM kernel service 106

## E

e\_assert\_wait kernel service 112

e\_block\_thread kernel service 113

e\_clear\_wait kernel service 114

e\_sleep kernel service 115

e\_sleep\_thread kernel service 117

e\_sleeppl kernel service 116

e\_wakeup kernel service 121

e\_wakeup\_one kernel service 121

e\_wakeup\_w\_result kernel service 121

e\_wakeup\_w\_sig kernel service 122

EEH Kernel Services

eeh\_broadcast 123

eeh\_clear 124

eeh\_disable\_slot 125

eeh\_enable\_dma 126

eeh\_enable\_pio 127

eeh\_enable\_slot 128

eeh\_init 129

eeh\_init\_multifunc 130

eeh\_read\_slot\_state 133

eeh\_reset\_slot 135

eeh\_slot\_error 136

eeh\_broadcast Kernel Service 123

eeh\_clear Kernel Service 124

eeh\_disable\_slot Kernel Service 125

eeh\_enable\_dma Kernel Service 126

eeh\_enable\_pio Kernel Service 127

eeh\_enable\_slot Kernel Service 128

eeh\_init Kernel Service 129

eeh\_init\_multifunc Kernel Service 130

eeh\_read\_slot\_state Kernel Service 133

eeh\_reset\_slot Kernel Service 135

eeh\_slot\_error Kernel Service 136

enable DMA

d\_map\_enable 81

End of Interrupt (EOI) kernel services

i\_eoi 210

enqueue kernel service 137

entry points

function pointers

obtaining 279

error logs

writing entries 139

error logs, writing entries 419

errresume kernel service 139

- errsave kernel service 139
- et\_post kernel service 119
- et\_wait kernel service 120
- event management
  - shared events
    - waiting for 115
- exception handlers
  - system-wide
    - deleting 510
    - systemwide 506
- exception information
  - retrieving 188
- exception management
  - contexts
    - saving 468
  - creating a process 54
  - execution flows
    - modifying 351
  - internationalized kernel message requests
    - submitting 386
  - locking 343
  - parent
    - setting to init process 469
  - putting process to sleep 476
  - sending a signal 406
  - states
    - saving 468
  - unmasked signals
    - determining if received 473
- exceptions 54
- execution flows
  - modifying 351
- execution states
  - saving 468
- ext parameter 614
- external storage
  - freeing 365

## F

- fetch\_and\_add kernel services 140
- fetch\_and\_and kernel service 141
- fetch\_and\_or kernel service 141
- fidtopv kernel service 142
- file attributes
  - getting 146
- file operation requirements 512
- file systems 149, 195
- file-mode creation mask 193
- files 159
  - access control lists
    - retrieving 662
    - setting 687
  - accessing blocks 691
  - attributes
    - getting 663
  - checking access permission 144
  - closing 145
  - creating 656
  - descriptor flags 192
  - descriptors 515, 516
  - determining if changed 562
  - hard links
    - requesting 667
  - interface to kernel services 511
  - mappings
    - validating 671

- files (*continued*)
  - opening 148, 151, 153, 156
  - opening for reading 675
  - opening for writing 675
  - pointers
    - retrieving 149
  - read subroutine 163
  - reading 163, 164, 166
  - readv subroutine 164
  - releasing portions of 658
  - renaming 682
  - size limit
    - retrieving 188
  - truncating 662
  - unlinking 681
  - unmapping 692
  - writing 166, 172
- find\_input\_type kernel service 143
- fp\_access kernel service 144
- fp\_close kernel service
  - GDLC 145
- fp\_close kernel service 145
  - device driver 620
- fp\_fstat kernel service 146
- fp\_fsync kernel service 147
- fp\_get\_path kernel service 150
- fp\_getdevno kernel service 148
- fp\_getea kernel service
  - opening regular files 148
- fp\_getf kernel service 149
- fp\_hold kernel service 151
- fp\_ioctl kernel service 151, 152
- fp\_ioctlx kernel service 153
- fp\_listea subroutine
  - support for 154
- fp\_lseek kernel service 155
- fp\_open kernel service
  - opening GDLC 157
  - opening regular files 156
- fp\_opendev kernel service 159
- fp\_poll kernel service 161
- fp\_read kernel service 163
- fp\_readv kernel service 164
- fp\_removeea kernel service
  - opening regular files 165
- fp\_rwuio kernel service 166
- fp\_select kernel service
  - cascaded support 167
  - invoking 168
  - notify routine and 168
  - returning from 169
- fp\_select kernel service notify routine 170
- fp\_setea kernel service
  - opening regular files 170
- fp\_statea kernel service 171
- fp\_write kernel service
  - data sent to DLC 173
  - open files 172
- fp\_writew kernel service 175
- free-pinned character buffers
  - sizing 410
- fskv\_reg kernel service 179
  - kv\_open() callout function 175
  - kv\_setattr() callout function 175
- fstatx subroutine
  - fp\_fstat kernel service 146
- fubyte kernel service 179

func subroutine 223  
fuword kernel service 180

## G

GDLC channels  
  disabling 145  
get\_pag Kernel Service 189  
get\_pag64 Kernel Service 189  
get\_umask kernel service 193  
getblk kernel service 181  
getc kernel service 182  
getcb kernel service 183  
getcbp kernel service 184  
getcf kernel service 185  
getcx kernel service 185  
geteblk kernel service 186  
geterror kernel service 187  
getexcept kernel service 188  
getfslimit kernel service 188  
getpid kernel service 190  
getppidx kernel service 190  
getuerror kernel service 191  
getuflags kernel service 192  
gfsadd kernel service 193  
gfsdel kernel service 195  
gn\_closecnt Subroutine 196  
gn\_common\_memcntl Subroutine 197  
gn\_mapcnt Subroutine 198  
gn\_opencnt Subroutine 199  
gn\_unmapcnt Subroutine 199  
groupmember Subroutine 200  
groupmember\_cr Subroutine 200

## H

heap\_create kernel service 201  
heap\_destroy kernel service 203  
heap\_modify kernel service 204  
heaps  
  initializing virtual memory 220  
hkeyset\_restore\_userkeys kernel service 205  
hkeyset\_update\_userkeys kernel service 206  
host names  
  obtaining 267  
hread\_set\_smtpriority system call 488

## I

i\_clear kernel service 207  
i\_disable kernel service 208  
i\_enable kernel service 209  
i\_eoi Kernel Service 210  
i\_init kernel service 217  
i\_mask kernel service 219  
i\_reset kernel service 237  
i\_sched kernel service 237  
i\_unmask kernel service 239  
I/O 182, 187, 207, 219  
  buffer cache  
    purging block from 422  
  buffers  
    freeing 426  
  character  
    retrieving 185

I/O (*continued*)  
  character buffer  
    waiting for free 583  
  character lists  
    using 619  
  characters  
    placing 423, 427  
  completion  
    waiting for 233  
  early power-off warning 218  
  free-pinned character buffers 410  
  freeing buffer lists 426  
  header memory buffers  
    allocating 370  
  interrupt handler  
    coding an 218  
  mbreq structures 356  
  mbuf chains  
    adjusting 372  
    appending 358  
    copying data from 363  
    freeing 366  
  mbuf clusters  
    allocating 360  
    allocating a page-sized 360  
  mbuf structures  
    allocating 359, 367, 368, 369, 370  
    attaching 369  
    clusters 371  
    converting pointers 376  
    creating 364  
    cross-memory descriptors 376  
    deregistering 364  
    freeing 365  
    initial requirements 373  
    pointers 375  
    removing 361  
    usage statistics 357  
  off-level processing  
    enabling 237  
  placing character buffers 424  
  placing characters 425  
I/O levels  
  waiting on 574  
identifiers  
  message queue 288  
idle to ready 221  
IDs  
  getting current process 190  
  getting parent 190  
if\_attach kernel service 213  
if\_detach kernel service 214  
if\_down kernel service 215  
if\_nostat kernel service 216  
ifa\_ifwithaddr kernel service 211  
ifa\_ifwithdstaddr kernel service 212  
ifa\_ifwithnet kernel service 213  
ifnet structures  
  address of 350  
ifunit kernel service 216  
init\_heap kernel service 220  
initp kernel service 221  
initp kernel service func subroutine 223  
input packets  
  building header for 439  
input types  
  adding new 8



- interface 213
- interface drivers
  - error handling 380
- interfaces
  - files 511
  - network
    - adding 214
- internationalized kernel message requests
  - submitting 386
- interrupt environment services
  - d\_cflush 62
  - getcx 185
  - if\_attach 214
  - net\_start\_done 383
  - tstart 498
- interrupt handlers 533
  - avoiding delays 237
  - coding 218
  - defining 217
  - removing 207
- interrupt priorities
  - disabling 208
  - enabling 209
- io\_map kernel service 223
- io\_map\_clear kernel service 224
- io\_map\_init kernel service 225
- io\_unmap kernel service 226
- iodone kernel service 227
- iodone routine
  - setting up 228
- iostadd kernel service 229
- iostdel kernel service 232
- iowait kernel service 233
- ip filtering hooks 234
- ip\_fltr\_in\_hook, ip\_fltr\_out, ipsec\_decap\_hook kernel service 234
- ipthreadsn 676, 678
- IS64U kernel service 239

## K

- k\_cpuextintr\_ctl kernel service 240
- kcap\_is\_set kernel service 241
- kcap\_is\_set\_cr kernel service 241
- kcid\_curproc kernel service 242
- kcred\_genpagvalue Kernel Service 243
- kcred\_getpag Kernel Service 245
- kcred\_getpag64 Kernel Service 245
- kcred\_getpagid kernel service 246
- kcred\_getpaginfo Kernel Service 246
- kcred\_getpagname kernel service 247
- kcred\_getppriv kernel service 248
- kcred\_setpag Kernel Service 251
- kcred\_setpag64 Kernel Service 251
- kcred\_setpagname kernel service 252
- kcred\_setppriv kernel service 253
- kern\_soaccept kernel service 254
- kern\_sobind kernel service 255
- kern\_soclose kernel service 256
- kern\_soconnect kernel service 257
- kern\_socreate kernel service 258
- kern\_sogetopt kernel service 259
- kern\_solisten kernel service 260
- kern\_soreceive kernel service 261
- kern\_soreserve kernel service 263
- kern\_sosend kernel service 264
- kern\_sosetopt kernel service 265

- kern\_soshutdown kernel service 266
- kernel buffers 615
- kernel memory 369
  - address ranges
    - pinning 354, 408, 605
    - releasing intersecting pages 565
    - setting storage protect key for 561
    - unpinning 355, 520, 606
  - address space
    - allocating 11
    - deallocating 12
    - deselecting 12
    - mapping 11, 21
    - obtaining handles 15
    - releasing 20
    - selecting 11
    - unmapping 12
  - bytes
    - retrieving 179
  - character data
    - copying into 49
  - characters
    - retrieving from buffers 531
    - writing to buffers 530
  - copying from 50
  - copying into 48
  - data
    - moving between VMO and buffer 551
    - retrieving a byte 179
    - retrieving a word 180
    - storing bytes 478
  - files
    - determining if changed 562
  - header memory buffers
    - allocating 370
  - heaps
    - initializing 220
  - I/O levels
    - waiting on 574
  - mbuf chains
    - adjusting 372
    - adjusting size of 356
    - appending 358
    - copying data from 363
    - freeing 366
    - reducing structures in 361
  - mbuf clusters
    - allocating 360
    - allocating a page-sized 360
  - mbuf structures
    - allocating 359, 367, 368, 370
    - attaching 369
    - clusters 371
    - converting addresses in 106
    - converting pointers 376
    - copying 362
    - creating 364
    - cross-memory descriptors 376
    - deregistering 364
    - freeing 365
    - initial requirements 373
    - pointers 375
    - removing 361
  - object modules
    - pinning 410
  - page ranges
    - initiating page-out 577

kernel memory (*continued*)

- page-out
  - determining I/O level 574
- page-ranges
  - initiating page-out 579
- pages
  - making without page-in 547
  - releasing several 566
- paging device tables
  - adding file system to 549
  - freeing entries in 577
- pin counts
  - decrementing 521
- storing words 480
- user buffer
  - preparing for access 599
- user-address space, 64-bit det 239
- virtual memory handles
  - constructing 543
- virtual memory manager 312
- virtual memory objects
  - creating 572
  - deleting 573
  - mapping to a region 536
- virtual memory resources
  - releasing 566
- words
  - retrieving 180

kernel messages
 

- printing to terminals 528

kernel object files
 

- loading 280
- unloading 283

kernel process state
 

- changing 221

kernel processes
 

- creation support 223

kernel service
 

- pin\_context\_stack 409
- unpin\_context\_stack 409

kernel services

- as\_att64 kernel service 11
- as\_det64 kernel service 12
- as\_geth kernel service 13
- as\_geth64 kernel service 14
- as\_getsrval64 kernel service 15
- as\_puth64 kernel service 20
- as\_seth64 kernel service 21
- bindprocessor 28
- compare\_and\_swap 45
- compare\_and\_swaplp 45
- disable\_lock 77
- e\_assert\_wait 112
- e\_block\_thread 113
- e\_clear\_wait 114
- e\_sleep\_thread 117
- e\_wakeup 121
- e\_wakeup\_one 121
- e\_wakeup\_w\_result 121
- e\_wakeup\_w\_sig 122
- et\_post 119
- et\_wait 120
- fetch\_and\_add 140
- fetch\_and\_addlp 140
- fetch\_and\_and 141
- fetch\_and\_or 141
- file interface to 511

kernel services (*continued*)

- IS64U 239
- kcred\_getpagid 246
- kcred\_getpagname 247
- kcred\_setpagname 252
- kthread\_kill 310
- kthread\_start 311
- limit\_sigs 336
- lock\_addr 345
- lock\_alloc 338
- lock\_clear\_recursive 339
- lock\_done 340
- lock\_free 341
- lock\_init 342
- lock\_islocked 342
- lock\_read 346
- lock\_read\_to\_write 346
- lock\_set\_recursive 348
- lock\_try\_read 346
- lock\_try\_read\_to\_write 346
- lock\_try\_write 348
- lock\_write 348
- lock\_write\_to\_read 349
- ltpin 354
- ltunpin 355
- rusage\_incr 461
- simple\_lock 474
- simple\_lock\_init 475
- simple\_lock\_try 474
- simple\_unlock 476
- thread\_create 485
- thread\_setsched 487
- thread\_terminate 490
- tstop 500
- tuning 501
- ufdgetf 515
- ufdhold 516
- ufdrele 516
- unlock\_enable 518
- user-mode exception handler for uexadd 507
- kgethostname kernel service 267
- kgetpname Kernel Service 268
- kgetrlimit64 kernel service 269
- kgetsystemcfg subroutine 271
- kgettckd kernel service 271
- kkey\_assign\_private kernel service 272
- kkeyset\_add\_key kernel service 273
- kkeyset\_add\_set kernel service 274
- kkeyset\_create kernel service 275
- kkeyset\_delete kernel service 275
- kkeyset\_remove\_key kernel service 276
- kkeyset\_remove\_set kernel service 277
- kkeyset\_to\_hkeyset kernel service 278
- klpar\_get\_info kernel service 278
- kmod\_entrypt kernel service 279
- kmod\_load kernel service 280
- kmod\_unload kernel service 283
- kmod\_util kernel service 284
- kmsgctl kernel service 287
- kmsgget kernel service 288
- kmsgsnd kernel service 291
- kmsrcv kernel service 290
- kprobe kernel service 419
- kra\_attachrset Subroutine 293
- kra\_creatp subroutine 295
- kra\_detachrset Subroutine 296
- kra\_getrset Subroutine 297

- krs\_alloc Subroutine 298
- krs\_free Subroutine 299
- krs\_get\_homesrad Subroutine 300
- krs\_getassociativity Subroutine 299
- krs\_getinfo Subroutine 301
- krs\_getpartition Subroutine 302
- krs\_getrad Subroutine 303
- krs\_init Subroutine 304
- krs\_numrads Subroutine 305
- krs\_op Subroutine 305
- krs\_setpartition Subroutine 307
- ksettickd kernel service 308
- ksettimer kernel service 309
- kthread\_kill kernel service 310
- kthread\_start kernel service 311
- kv\_open() callout function 175
- kv\_setattr() callout function 175
- kvm\_pattr Kernel Service 554
- kvmgetinfo kernel service 312
- kwpar\_checkpoint\_status kernel service 314
- kwpar\_err kernel service 315
- kwpar\_getname Kernel Service 316
- kwpar\_getrootpath Kernel Service 317
- kwpar\_isappwpar Kernel Service 318
- kwpar\_r2vmap\_devno kernel service 318
- kwpar\_r2vmap\_pid kernel service 319
- kwpar\_r2vmap\_tid kernel service 320
- kwpar\_regdevno kernel service 321
- kwpar\_reghook kernel service 322
- kwpar\_unregdevno kernel service 324
- kwpar\_unreghook kernel service 325
- kwpar\_v2rmap\_devno kernel service 326
- kwpar\_v2rmap\_pid kernel service 327
- kwpar\_v2rmap\_tid kernel service 328

## L

- ldata\_alloc Kernel Service 329
- ldata\_create Kernel Service 329
- ldata\_destroy Kernel Service 331
- ldata\_free Kernel Service 332
- ldata\_grow Kernel Service 332
- ldmp\_bufest kernel service 333
- ldmp\_errstr kernel service 333
- ldmp\_freeparms kernel service 335
- ldmp\_setupparms kernel service 335
- ldmp\_timeleft kernel service 333
- ldmp\_xmalloc kernel service 333
- ldmp\_xmfree kernel service 333
- limit\_sigs kernel service 336
- livedump kernel service 337
- lock\_addr kernel service 345
- lock\_alloc kernel service 338
- lock\_clear\_recursive kernel service 339
- lock\_done kernel service 340
- lock\_free kernel service 341
- lock\_init kernel service 342
- lock\_islocked kernel service 342
- lock\_read kernel service 346
- lock\_read\_to\_write kernel service 346
- lock\_set\_recursive kernel service 348
- lock\_try\_read kernel service 346
- lock\_try\_read\_to\_write kernel service 346
- lock\_try\_write kernel service 348
- lock\_write kernel service 348
- lock\_write\_to\_read kernel service 349
- locking 43

- lockl kernel service 343
- logical file system
  - channel numbers
    - finding 148
  - device numbers
    - finding 148
  - file attributes
    - getting 146
  - file descriptors
    - status of 161
  - file pointers
    - retrieving 149
    - status of 161
  - files
    - checking access permissions 144
    - closing 145
    - opening 148, 151, 153, 156
    - reading 164, 166
    - writing 166, 172, 175
  - fp\_listea subroutine
    - support for 154
  - fp\_removeea subroutine
    - support for 165
  - message queues
    - status of 161
  - notify routine
    - registering 170
  - offsets
    - changing 155
  - open subroutine
    - support for 148, 156, 171, 172
  - poll request 167
  - read subroutine
    - interface to 163
  - readv subroutine
    - interface to 164
  - select operation 167
  - special files
    - opening 159
  - use count
    - incrementing 151
  - write subroutine 172
  - writew subroutine
    - interface to 175
- loifp kernel service 350
- longjmpx kernel service 351
- lookupname kernel service 351
- lookupname\_cur kernel service 351
- lookupvp kernel service 351
- looutput kernel service 353
- ltpin kernel service 354
- ltunpin kernel service 355

## M

- m\_adj kernel service 356
- m\_cat kernel service 358
- m\_clattach kernel service 359
- m\_clget macro 360
- m\_clgetm kernel service 360
- m\_collapse kernel service 361
- m\_copy macro 362
- m\_copydata kernel service 363
- m\_copym kernel service 364
- m\_dereg kernel service 364
- m\_freem kernel service 366
- m\_get kernel service 367

- m\_getclr kernel service 368
- m\_getclust macro 368
- m\_getclustm kernel service 369
- m\_gethdr kernel service 370
- M\_HASCL kernel service 371
- m\_pullup kernel service 372
- m\_reg kernel service 373
- M\_XMEMD macro 376
- macros
  - add\_netopt 10
  - del\_netopt 67
  - DTOM 106
  - m\_clget 360
  - m\_getclust 368
  - M\_HASCL 371
  - MTOCL 375
  - MTOD 376
- maps DMA master devices
  - d\_map\_page 86
- mbreq structure
  - format of 356
- mbuf chains
  - adjusting 372
  - adjusting size of 356
  - appending 358
  - copying 363
  - freeing 366
  - removing structures from 361
- mbuf clusters
  - allocating 360
  - allocating a page-sized 360
  - page-sized
    - attaching 369
- mbuf structures
  - address to header 106
  - allocating 359, 367, 368, 370
  - attaching a cluster 369
  - clusters
    - determining presence of 371
  - converting pointers 376
  - copying 362, 364
  - cross-memory descriptors
    - obtaining address of 376
  - deregistering 364
  - freeing 365
  - initial requirements 373
  - mbreq structure 356
  - mbstat structure 357
  - pointers
    - converting 375
  - registration information 356
  - removing 361
  - usage statistics 357
- memory
  - allocating 598
  - buffers (device drivers) 639
  - freeing 611
  - pages
    - preparing for DMA 602, 603
    - processing after DMA I/O 602, 603
  - performing a cross-memory move 608, 609
  - rmfree 450
  - uio structures 639
  - user buffer
    - detaching from 601
- memory allocation
  - rmalloc 449

- memory manager
  - kvmgetinfo 312
- memory mapped I/O
  - rmmmap\_create 450
  - rmmmap\_remove 454
- message queues
  - control operations
    - providing 287
  - identifiers
    - obtaining 288
  - messages
    - reading 290
    - sending 291
- MTOCL macro 375
- MTOD macro 376
- multiplexed device driver
  - allocating 627
  - deallocating 627
- mycpu kernel service 377

## N

- nameToXfid() kernel service 377
- net\_attach kernel service 378
- net\_detach kernel service 379
- net\_error kernel service 380
- net\_sleep kernel service 381
- net\_start kernel service 382
- net\_start\_done kernel service 383
- net\_wakeup kernel service 383
- net\_xmit kernel service 384
- net\_xmit\_trace kernel service 385
- network
  - ctlinput function
    - invoking 404
  - current host name 267
  - demuxers
    - adding 388
    - deleting 395
    - disabling 396
    - enabling 389
  - destination addresses
    - locating 212
  - device drivers
    - allocating 393
    - relinquishing 398
  - device handlers
    - closing 379
    - ending a start 383
    - opening 378
    - starting ID on 382
  - devices
    - attaching 394
    - detaching 398
  - ID
    - ending a start 383
  - ifnet structures
    - address of 350
  - input packets
    - building header for 439
  - interface
    - adding 213
  - interface drivers
    - error handling 380
  - putting caller to sleep 381
  - raw protocols
    - implementing user requests for 439

- network (*continued*)
  - raw\_header structures
    - building 439
  - receive filters
    - adding 389
    - deleting 396
  - routes
    - allocating 455
  - routing table entries
    - changing 459, 460
    - creating 457
    - forcing through gateway 458
    - freeing 456
  - software interrupt service routines
    - invoking 462
    - scheduling 462
  - start operation
    - ending 383
  - status filters
    - adding 392
    - deleting 397
  - transmit packets
    - tracing 385
  - waking sleeping processes 383
- network address families
  - adding 8
  - deleting 64
  - searching for 405
- network device handlers
  - transmitting packets 384
- network input types
  - adding 8
  - deleting 65
- network interfaces
  - deleting 214
  - locating 211, 213
  - marking as down 215
- pointers
  - obtaining 216
- software loopback
  - obtaining address 350
  - sending data through 353
  - zeroing statistic elements 216
- network option structures
  - adding 10
  - deleting 67
- network packet types
  - finding 143
- network software interrupt service
  - adding 10
  - deleting 66
- NLuprint kernel service 386
- notify routine
  - registering 170
    - from fp\_select kernel service 168
- ns\_add\_demux network service 388
- ns\_add\_filter network service 389
- ns\_add\_status network service 392
- ns\_alloc network service 393
- ns\_attach network service 394
- ns\_del\_demux network service 395
- ns\_del\_filter network service 396
- ns\_del\_status network service 397
- ns\_detach network service 398
- ns\_free network service 398

## O

- object modules
  - pinning 410
- off-level processing 237
- offset
  - changing 155
- open subroutine
  - support for 148, 156

## P

- packet types
  - finding 143
- packets
  - transmitting 384
- page-out
  - determining I/O level 574
- page-ranges
  - initiating page-out 577
- pages
  - making without page-in 547
  - releasing several 566
- paging device tables
  - adding file system to 549
  - freeing entries in 577
- panic kernel service 403
- PCI bus slot configuration registers 403
- pci\_cfgwr kernel service 403
- pfctlinput kernel service 404
- pfindproto kernel service 405
- pgsignal kernel service 406
- pidsig kernel service 407
- pin counts
  - decrementing 521
- pin kernel service 408
- pin\_context\_stack kernel service 409
- pincf kernel service 410
- pincode kernel service 410
- pio\_assist kernel service 411
- pipes
  - select request on 167
- poll request
  - registering asynchronous 465
  - support for 463
- power-off warnings
  - registering early 218
- privileges
  - checking effective 479
- probe kernel service 419
- process 54
- process environment services
  - d\_cflush 62
  - ddread entry point 631
  - getcx 185
  - i\_disable 208
  - if\_attach 214
  - iostdel 232
  - net\_attach 378
  - net\_start\_done 383
  - tstart 498
- process management
  - blocking a process 509
  - calling process IDs 190
  - checking effective privileges 479
  - clearing blocked processes 510

- process management (*continued*)
  - contexts
    - removing 43
    - saving 468
  - creating a process 54
  - execution flows
    - modifying 351
  - forcing a wait 115
  - idle to ready
    - transition of 221
  - internationalized kernel message requests
    - submitting 386
  - locking 343
  - parent
    - setting to init process 469
  - parent process IDs
    - getting 190
  - process initialization routine
    - directing 223
  - process state-change notification routine 417
  - putting process to sleep 476
  - shared events
    - waiting for 116
  - signals
    - sending 406, 407
  - signals, sending 407
  - state transition notification 414
  - state-change notification routine
    - deleting 419
  - states
    - saving 468
  - suspending processing 64
  - unlocking
    - conventional processes 519
  - unmasked signals
    - determining if received 473
  - wait
    - for shared event 116
  - waking up processes 463
- process state-change notification routine 414
- processor cache
  - flushing 537
- proch structure 418
- proch\_unreg kernel service 417
- prochadd kernel service 417
- prochdel kernel service 419
- programmed I/O
  - exceptions caused by 411
- purblk kernel service 422
- putc kernel service 423
- putcb kernel service 424
- putcbp kernel service 425
- putcf kernel service 426
- putcfl kernel service 426
- putcx kernel service 427

## Q

- query\_proc\_info kernel service 428
- queue elements
  - checking validity 41
  - cleanup 37
  - placing into queue 137
  - waiting for 583
- queue management routines
  - attach-device 22
  - cancel-queue-element 37

- queue management routines (*continued*)
  - detach-device 68
  - parameter checking 41

## R

- RAS kernel services
  - error logs
    - writing entries in 139
  - master dump table
    - deleting entry from 100
- RAS services
  - system crash
    - performing system dump of 403
  - trace events
    - recording 493, 494
- RAS\_BLOCK\_NULL Exported Data Structure 429
- ras\_control Exported Kernel Service 429
- ras\_customize Exported Kernel Service 430
- ras\_path\_control Exported Kernel Service 431
- ras\_register Exported Kernel Service 432
- ras\_ret\_query\_parms kernel service 434
- ras\_unregister Exported Kernel Service 432
- raschk\_eaddr\_hkeyset kernel service 435
- raschk\_eaddr\_kkey kernel service 436
- raschk\_stktrace kernel service 437
- raw protocols
  - implementing user requests for 439
- raw\_header structures
  - building 439
- raw\_input kernel service 439
- raw\_usrreq kernel service 439
- rawinch field 230
- read subroutine
  - interface to 163
- read-ahead block
  - starting I/O on 31
- readv subroutine
  - interface to 164
- ready to idle 221
- reconfig\_complete kernel service 441
- reconfig\_register kernel service 441
- reconfig\_register\_ext kernel service 441
- reconfig\_unregister kernel service 441
- record locking 43
- record locks
  - controlling 668
- refmon kernel service 446
- regions
  - unmapping virtual memory 537
- Resource Set APIs
  - kra\_attachrset 293
  - kra\_creatp 295
  - kra\_detachrset 296
  - kra\_getrset 297
  - krs\_alloc 298
  - krs\_free 299
  - krs\_get\_homesrad 300
  - krs\_getassociativity 299
  - krs\_getinfo 301
  - krs\_getpartition 302
  - krs\_getrad 303
  - krs\_init 304
  - krs\_numrads 305
  - krs\_op 305
  - krs\_setpartition 307

- resources
  - virtual file system
    - releasing 535
- rmalloc kernel service 449
- rmfree kernel service 450
- rmmmap\_create kernel service 450
- rmmmap\_remove kernel service 454
- routes
  - allocating 455
- routing table entries
  - changing 459, 460
  - creating 457
  - forcing through gateway 458
  - freeing 456
- rtalloc kernel service 455
- rtfree kernel service 456
- rtinit kernel service 457
- rtredirect kernel service 458
- rtrequest kernel service 459, 460
- rusage\_incr kernel service 461

## S

- schednetisr kernel service 462
- scheduling functions 490
- security subroutines
  - kcred\_genpagvalue 243
- select request
  - registering asynchronous 465
  - support for 463
- selnotify kernel service 463
- selreg kernel service 465
- set\_pag Kernel Service 467
- set\_pag64 Kernel Service 467
- setioctlrv Subroutine 467
- setjmpx kernel service 468
- setpinit kernel service 469
- setuerror kernel service 469
- setufdflags kernel service 192
- shared events
  - waiting for 116
- shared memory
  - controlling access to 343
- shared object modules
  - symbol resolution 282
- shutdown kernel services
  - shutdown\_notify\_reg 470
  - shutdown\_notify\_unreg 472
- shutdown\_notify\_reg kernel service 470
- shutdown\_notify\_unreg kernel service 472
- sig\_chk kernel service 473
- signals
  - sending 406
- simple\_lock kernel service 474
- simple\_lock\_init kernel service 475
- simple\_lock\_try kernel service 474
- simple\_unlock kernel service 476
- sleep kernel service 476
- sockets
  - select request on 167
- software interrupt service routines
  - invoking 462
  - scheduling 462
- software loopback interfaces
  - obtaining address of 350
  - sending data through 353
- software-interrupt level 10

- special files
  - creating 675
  - opening 159
  - requesting I/O control operations 666
- standard parameters
  - device driver 614
- statistics structures
  - registering 229
  - removal 232
- strategy routine
  - calling 526
- subyte kernel service 478
- suser kernel service 479
- suword kernel service 480
- switch table 75
- symbol binding support 282
- symbol resolution and shared object modules 282
- symbolic links
  - reading contents of 680
- synchronization functions
  - providing 271
- system call events
  - auditing 24
- system calls
  - \_\_pag\_getid 399
  - \_\_pag\_getname 400
  - \_\_pag\_getvalue 400
  - \_\_pag\_setname 401
  - \_\_pag\_setvalue 402
  - thread\_set\_smtpriority 488
- system dump kernel services
  - dmp\_add 92
  - dmp\_ctl 96
- system dumps
  - adding and removing master dump table entries 96
  - adding to master dump table 92
  - performing 403
  - specifying contents 92
- systemwide time
  - setting 309

## T

- talloc kernel service 484
- TE\_verify\_reg kernel service 481
- TE\_verify\_unreg kernel service 483
- tfree kernel service 484
- thread\_create kernel service 485
- thread\_read\_smtpriority system calls 488
- thread\_self subroutine 486
- thread\_setsched kernel service 487
- thread\_terminate kernel service 490
- time
  - allocating time request blocks 484
  - callout table entries
    - registering changes in 492
  - canceling pending timer requests 523
  - current
    - reading 60
  - scheduling functions 490
  - submitting timer request 498
  - suspending processing 64
  - synchronization functions
    - providing 271
  - systemwide
    - setting 309

- time (*continued*)
  - time request blocks
    - deallocating 484
  - time-adjustment value 271
  - updating 308
  - watchdog timers
    - registering 587
    - removing 586, 588
    - stopping 590
- timeout kernel service 490
- timeoutcf kernel subroutine 492
- timer
  - watchdog timers
    - starting 589
- trace events
  - recording 493, 494, 495
- transfer requests
  - tailoring 527
- transmit packets
  - tracing 385
- trc\_ishookon Exported Kernel Service 493
- trcgenk kernel service 493
- trcgenkt kernel service
  - DLC 495
  - recording for a generic trace channel 494
- tstart kernel service 498
- tstop kernel service 500
- tty device driver support 230
- ttystat structure 229
- tuning kernel service 501

## U

- ue\_proc\_check kernel service 504
- ue\_proc\_register subroutine 505
- ue\_proc\_unregister subroutine 506
- uexadd kernel service
  - adding an exception handler 506
- uexblock kernel service 509
- uexclear kernel service 510
- uexdel kernel service 510
- ufdcreate kernel service 511
- ufdgetf kernel service 515
- ufdhold kernel service 516
- ufdrele kernel service 516
- uio structures 384, 639
- uiomove kernel service 517
- unlock\_enable kernel service 518
- unlocking conventional processes 519
- unlockl kernel service 519
- unpin kernel service 520
- unpin\_context\_stack kernel service 409
- unpincode kernel service 521
- untimeout kernel service 523
- uphysio kernel mincnt service 527
- uphysio kernel service
  - described 524
  - error detection by 526
  - mincnt routine 527
- uprintf kernel service 528
- uprintf structure 386
- ureadc kernel service 530
- use count
  - incrementing 151
- user buffer
  - detaching from 601
  - preparing for access 599

- user-address space 239
- user-mode exception handler for uexadd kernel service 507
- ut\_error field
  - retrieving 191
- ut\_error fields
  - setting 469
- uwritec kernel service 531

## V

- v-node operations 663, 666, 670, 674, 683, 684
  - retrieving 351
- v-nodes 663
  - allocating 580
  - closing associated files 655
  - count
    - incrementing 665
  - file identifier conversion to 652
  - file identifiers
    - building 659
  - finding by name 670
  - freeing 580
  - modifications
    - flushing to storage 660
  - obtaining root 648
  - polling 686
  - releasing references 681
  - validating access to 653
- validate\_pag Kernel Service 532
- validate\_pag64 Kernel Service 532
- vec\_clear kernel service 533
- vec\_init kernel service 534
- VFS 663
  - access control lists
    - retrieving 662
  - allocating virtual nodes 580
  - building file identifiers 659
  - changes
    - writing to storage 650
  - checking record locks 668
  - control operations
    - implementing 644
  - creating directories 674
  - creating special files 675
  - file attributes
    - getting 663
  - file system types
    - adding 193
    - removing 195
  - files
    - accessing blocks 691
    - converting identifiers 652
    - creating 656
    - hard links 667
    - opening 675
    - releasing portions of 658
    - renaming 682
    - requesting I/O 676
    - setting access control 687
    - setting attributes 688
    - truncating 662
    - validating mapping requests 671
  - finding v-nodes by name 670
  - flushing v-node modifications 660
  - freeing virtual nodes 580
  - incrementing v-node counts 665
  - initializing 645



- VFS (*continued*)
  - mounting 646
  - nodes
    - pointer to root 648
    - retrieving 351
  - polling v-nodes 686
  - querying record locks 668
  - reading directory entries 678
  - releasing v-node references 681
  - removing directories 684
  - renaming directories 682
  - resources
    - releasing 535
  - revoking access 683
  - searching 648
  - setting record locks 668
  - special files
    - I/O control operations on 666
  - statistics
    - obtaining 649
  - structures, holding and releasing 645
  - unmounting 651
- VFS operations
  - vfs\_cntl 644
  - vfs\_hold 645
  - vfs\_init 645
  - vfs\_mount 646
  - vfs\_root 648
  - vfs\_search 648
  - vfs\_statfs 649
  - vfs\_sync 650
  - vfs\_umount 651
  - vfs\_unhold 645
  - vfs\_vget 652
  - vnop\_access 653
  - vnop\_close 655
  - vnop\_create 656
  - vnop\_fclear 658
  - vnop\_fid 659
  - vnop\_fsync 660
  - vnop\_fsync\_range 660
  - vnop\_ftruncate 662
  - vnop\_getacl 662
  - vnop\_hold 665
  - vnop\_link 667
  - vnop\_lockctl 668
  - vnop\_mknod 675
  - vnop\_open 675
  - vnop\_rdw 676
  - vnop\_rdw\_attr 676
  - vnop\_readdir 678
  - vnop\_readlink 680
  - vnop\_remove 681
  - vnop\_rename 682
  - vnop\_select 686
  - vnop\_setacl 687
  - vnop\_setattr 688
  - vnop\_strategy 691
  - vnop\_symlink 692
  - vnop\_unmap 692
- vfs\_aclxctl entry point 642
- vfsrele kernel service 535
- virtual file system 193, 662
- virtual interrupt handlers
  - defining 534
  - removing 533

- virtual memory
  - regions
    - unmapping 537
- virtual memory handles
  - constructing 543
- virtual memory objects
  - creating 572
  - deleting 573
  - managing addresses 11
  - mapping 21
  - mapping to a region 536
  - obtaining handles 13, 14, 15
  - page-out for range in 579
  - releasing 20
  - unmapping 12
- virtual memory resources
  - releasing 565
- vm\_att kernel service 536
- vm\_cflush kernel service 537
- vm\_det kernel service 537
- vm\_flushp kernel service 538
- vm\_guatt Kernel Service 541
- vm\_gudet Kernel Service 542
- vm\_handle kernel service 543
- vm\_invalidatep kernel service 544
- vm\_ioaccessp kernel service 545
- vm\_makep kernel service 547
- vm\_mount kernel service 549
- vm\_mounte kernel service 550
- vm\_mvc kernel service 553
- vm\_patr System Call 554
- vm\_protect\_kkey kernel service 560
- vm\_protectp kernel service 561
- vm\_qmodify kernel service 562
- vm\_qpages kernel service 563
- vm\_readp kernel service 564
- vm\_release kernel service 565
- vm\_releasep kernel service 566
- vm\_segmap kernel service 567
- vm\_setdevid kernel service 568
- vm\_setseg\_kkey kernel service 570
- vm\_thrpio\_pop kernel service 570
- vm\_thrpio\_push kernel service 571
- vm\_umount kernel service 577
- vm\_write kernel service 577
- vm\_writep kernel service 579
- vms\_create kernel service 572
- vms\_delete kernel service 573
- vms\_iowait kernel service 574
- vms\_iowaitf kernel service 574
- vn\_free kernel service 580
- vn\_get kernel service 580
- vnop\_getxcl entry point 664
- vnop\_ioctl entry point 666
- vnop\_memcntl Entry Point 673
- vnop\_seek Entry Point 685
- vnop\_setxcl entry point 689
- vnop\_symlink entry point 692
- vsx\_disable kernel service 581
- vsx\_enable kernel service 582

## W

- w\_clear kernel service 586
- w\_init kernel service 587
- w\_setattr kernel service 588
- w\_start kernel service 589

- w\_stop kernel service 590
- wait channels
  - putting caller to sleep on 381
- waitcfree kernel service 583
- waiting for free buffer 583
- waitq kernel service 583
- waking sleeping processes 383
- watchdog timers
  - registering 587
  - removing 586, 588
  - starting 589
  - stopping 590
- words
  - retrieving 180
  - storing in kernel memory 480
- write subroutine
  - interface to 172
- writev subroutine
  - interface to 175

## X

- xfidToName() kernel service 591
- xlate\_create kernel service 591
- xlate\_pin kernel service 593
- xlate\_remove kernel service 594
- xlate\_unpin kernel service 595
- xm\_det kernel service 596
- xm\_mapin 596
- xm\_maxmap Kernel Service 597
- xmalloc kernel service
  - described 598
- xmattach kernel service 599
- xmdetach kernel service 601
- xmemdma kernel service 602
- xmemdma64 kernel service 603
- xmemin kernel service 608
- xmemout kernel service 609
- xmempin kernel service 605
- xmempsize Kernel Service 610
- xmemunpin kernel service 606
- xmemzero kernel service 607
- xmfree kernel service 611
- xmgethkeyset kernel service 612
- xmsethkeyset kernel service 613





Printed in USA