

RED HAT CEPH STORAGE HARDWARE CONFIGURATION GUIDE

Designing scalable workload-optimized Ceph clusters



Ceph has been developed to deliver object, file, and block storage in one self-managing, self-healing platform with no single point of failure.

Red Hat® Ceph Storage offers multi-petabyte software-defined storage for the enterprise, across a range of industry-standard hardware.

With proper configuration, Red Hat Ceph Storage clusters can be designed for IOPS-optimized, throughput-optimized, or cost/capacity-optimized workloads.

EXECUTIVE SUMMARY

Ceph users frequently request simple, optimized cluster configurations for different workload types. Common requests are for throughput-optimized and capacity-optimized workloads, but IOPS-intensive workloads on Ceph are also emerging. Based on extensive testing by Red Hat with a variety of hardware providers, this document provides general performance, capacity, and sizing guidance. The companion “Red Hat Ceph Storage Hardware Selection Guide” provides sample hardware configurations sized to specific workloads. In-depth performance and sizing guides for Red Hat® Ceph Storage are also available for key hardware vendors.

TABLE OF CONTENTS

1 INTRODUCTION	2
2 CEPH ARCHITECTURE OVERVIEW	3
3 CLUSTER CONFIGURATION GUIDANCE	4
3.1 Qualifying the need for scale-out storage	4
3.2 Identifying target workload I/O profiles	5
3.3 Choosing a storage access method	6
3.4 Identifying capacity needs	7
3.5 Determining fault domain risk tolerance	8
3.6 Selecting a data protection method	9
4 HARDWARE CONFIGURATION GUIDELINES	10
4.1 Monitor nodes	10
4.2 OSD hosts	11
4.3 Broad OSD host configuration trends	14
5 CONCLUSION	15



facebook.com/redhatinc
@redhatnews
linkedin.com/company/red-hat

INTRODUCTION

Storage infrastructure is undergoing tremendous change, particularly as organizations deploy storage to support big data and private clouds. Traditional scale-up arrays are limited in scalability, and complexity can compromise cost-effectiveness. In contrast, software-defined storage infrastructure based on clustered storage servers has emerged as a way to deploy cost-effective and manageable storage at scale, with Ceph among the leading solutions. In fact, cloud storage companies are already using Ceph at near exabyte scale, with expected continual growth. For example, Yahoo estimates that their Ceph-based cloud object store will grow 20-25% annually.¹

Red Hat Ceph Storage significantly lowers the cost of storing enterprise data and helps organizations manage exponential data growth. The software is a robust, petabyte-scale storage platform for those deploying public or private clouds. As a modern storage system for cloud deployments, Red Hat Ceph Storage offers mature interfaces for enterprise block and object storage, making it well suited for active archive, rich media, and cloud infrastructure workloads like OpenStack®.² Delivered in a unified self-healing and self-managing platform with no single point of failure, Red Hat Ceph Storage handles data management so businesses can focus on improving application availability, with properties that include:

- Scaling to exabytes
- No single point of failure in the cluster
- Lower capital expenses (CapEx) by running on commodity server hardware
- Lower operational expenses (OpEx) by self-managing and self-healing

Many organizations are trying to understand how to configure hardware for optimized Ceph clusters that meet their unique needs. Red Hat Ceph Storage is able to run on myriad diverse industry-standard hardware configurations, but designing a successful Ceph cluster requires careful analysis of issues related to application, capacity, workload. The ability to address dramatically different kinds of I/O workloads within a single Ceph cluster makes understanding these issues paramount to a successful deployment.

After extensive performance and server scalability evaluation and testing with many vendors, Red Hat has developed a proven methodology that helps ask and answer key questions that lead to properly sized and configured scale-out storage clusters based on Red Hat Ceph Storage. Described in greater detail in this guide, the methodology includes:

- Qualifying the need for scale-out storage
- Identifying target workload I/O profiles
- Choosing a storage access method
- Identifying capacity
- Determining fault domain risk tolerance
- Selecting a data protection method

¹ <http://yahooeng.tumblr.com/post/116391291701/yahoo-cloud-object-store-object-storage-at>

² Ceph is and has been the leading storage for OpenStack according to several semi-annual OpenStack user surveys.

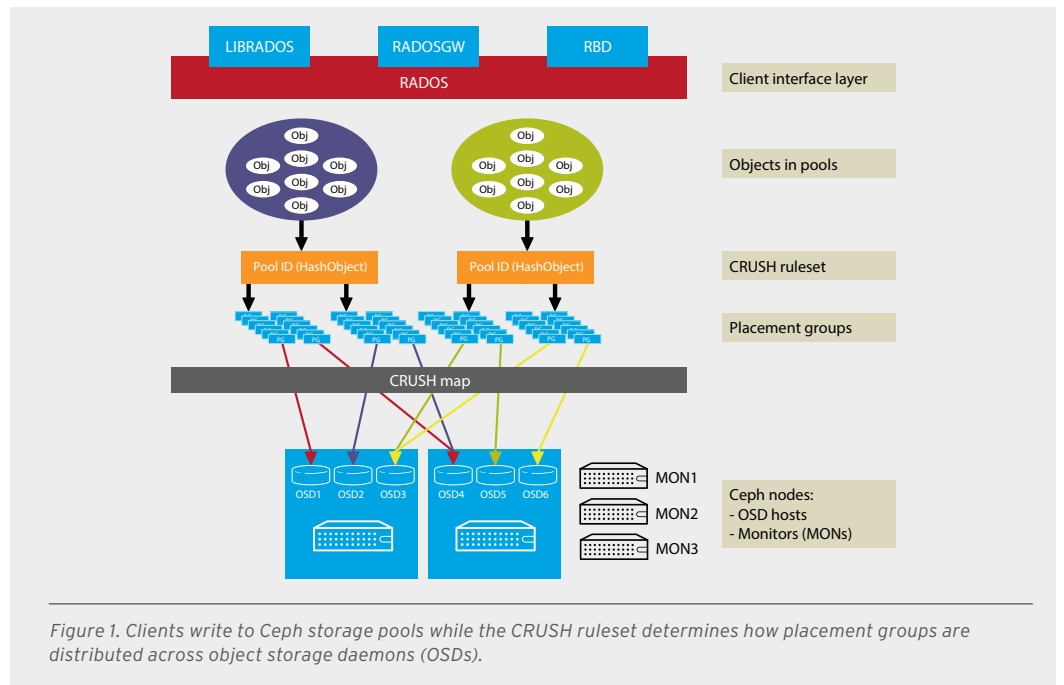
CEPH ARCHITECTURE OVERVIEW

A Ceph storage cluster is built from large numbers of Ceph nodes for scalability, fault-tolerance, and performance. Each node is based on commodity hardware and uses intelligent Ceph daemons that communicate with each other to:

- Store and retrieve data
- Replicate data
- Monitor and report on cluster health
- Redistribute data dynamically (remap and backfill)
- Ensure data integrity (scrubbing)
- Detect and recover from faults and failures

To the Ceph client interface that reads and writes data, a Ceph storage cluster looks like a simple pool where data is stored. However, the storage cluster performs many complex operations in a manner that is completely transparent to the client interface. Ceph clients and Ceph object storage daemons (Ceph OSD daemons, or OSDs) both use the CRUSH (controlled replication under scalable hashing) algorithm for storage and retrieval of objects.

For a Ceph client, the storage cluster is very simple. When a Ceph client reads or writes data (referred to as an I/O context), it connects to a logical storage pool in the Ceph cluster. Figure 1 illustrates the overall Ceph architecture, with concepts that are described in the sections that follow.



- **Pools.** A Ceph storage cluster stores data objects in logical dynamic partitions called pools. Pools can be created for particular data types, such as for block devices, object gateways, or simply to separate user groups. The Ceph pool configuration dictates the number of object replicas and the number of placement groups (PGs) in the pool. Ceph storage pools can be either replicated or erasure coded, as appropriate for the application and cost model. Additionally, pools can “take root” at any position in the CRUSH hierarchy, allowing placement on groups of servers with differing performance characteristics—allowing storage to be optimized for different workloads.
- **Placement groups.** Ceph maps objects to placement groups (PGs). PGs are shards or fragments of a logical object pool that are composed of a group of Ceph OSD daemons that are in a peering relationship. Placement groups provide a means of creating replication or erasure coding groups of coarser granularity than on a per object basis. A larger number of placement groups (e.g., 200 per OSD or more) leads to better balancing.
- **CRUSH ruleset.** The CRUSH algorithm provides controlled, scalable, and declustered placement of replicated or erasure-coded data within Ceph and determines how to store and retrieve data by computing data storage locations. CRUSH empowers Ceph clients to communicate with OSDs directly, rather than through a centralized server or broker. By determining a method of storing and retrieving data by algorithm, Ceph avoids a single point of failure, a performance bottleneck, and a physical limit to scalability.
- **Ceph monitors (MONs).** Before Ceph clients can read or write data, they must contact a Ceph MON to obtain the current cluster map. A Ceph storage cluster can operate with a single monitor, but this introduces a single point of failure. For added reliability and fault tolerance, Ceph supports an odd number of monitors in a quorum (typically three or five for small to mid-sized clusters). Consensus among various monitor instances ensures consistent knowledge about the state of the cluster.
- **Ceph OSD daemons.** In a Ceph cluster, Ceph OSD daemons store data and handle data replication, recovery, backfilling, and rebalancing. They also provide some cluster state information to Ceph monitors by checking other Ceph OSD daemons with a heartbeat mechanism. A Ceph storage cluster configured to keep three replicas of every object requires a minimum of three Ceph OSD daemons, two of which need to be operational to successfully process write requests. Ceph OSD daemons roughly correspond to a file system on a physical hard disk drive.

CLUSTER CONFIGURATION GUIDANCE

Despite the flexibility of Ceph, no one cluster or pool configuration fits all applications or situations. Instead, successful configuration of a Ceph cluster requires answering key questions about how the cluster will be used and the applications it will serve.

QUALIFYING THE NEED FOR SCALE-OUT STORAGE

Not every storage situation calls for scale-out storage. When requirements include several of the following needs, scale-out storage is likely the best solution.

- **Dynamic storage provisioning.** By dynamically provisioning capacity from a pool of storage, organizations are typically building a private storage cloud, emulating services such as Amazon Simple Storage Service (S3) for object storage or Amazon Elastic Block Store (EBS).
- **Standard storage servers.** Scale-out storage employs storage clusters built from industry-standard x86 servers rather than proprietary storage appliances, allowing incremental growth of storage capacity and/or performance without forklift appliance upgrades.

- **Unified name spaces.** Scale-out storage allows pooling storage across tens, hundreds, or even thousands of storage servers in one or more unified namespaces.
- **High data availability.** Scale-out storage provides high-availability of data across what would otherwise be “server storage islands” within the storage cluster.
- **Independent scalability of performance and capacity.** Unlike typical scale-up NAS and SAN devices which frequently run out of performance before running out of capacity, scale-out storage allows organizations to add storage performance or capacity incrementally by independently adding more storage servers or disks as required.

IDENTIFYING TARGET WORKLOAD I/O PROFILES

Accommodating the target workload I/O profile is perhaps the most crucial design consideration. As a first approximation, organizations need to understand if they are simply deploying low-cost archive storage or if their storage needs to meet specific performance requirements. For performance-oriented Ceph clusters, IOPS, throughput, and latency requirements must be clearly defined. On the other hand, if the lowest cost per terabyte is the overriding need, a Ceph cluster architecture can be designed at dramatically lower costs. For example, Ceph object archives with erasure-coded pools and without dedicated SSD write journals can be dramatically lower in cost than Ceph block devices on 3x-replicated pools with dedicated flash write journals.

If needs are more performance-oriented, IOPS and throughput are often taken into consideration. Historically, Ceph has performed very well with high-throughput workloads, and has been widely deployed for these use cases. Use cases are frequently characterized by large-block, asynchronous, sequential I/O (e.g., digital media performance nodes). In contrast, high IOPS workloads are frequently characterized by small-block synchronous random I/O (e.g., 4 KB random I/O). At present, the use of Ceph for high IOPS open source database workloads is emerging (e.g., MySQL, MariaDB, and PostgreSQL). Moreover, when Ceph is deployed as Cinder block storage for OpenStack virtual machine (VM) instances, it typically serves a mix of IOPS- and throughput-intensive I/O patterns.

Additionally, understanding the workload read/write mix can affect architecture design decisions. For example, erasure-coded pools can perform better than replicated pools for sequential writes, and worse than replicated pools for sequential reads. As a result, a write-mostly object archive workload (like video surveillance archival) may perform similarly between erasure-coded pools and replicated pools, with erasure-coded pools being significantly less expensive.

To simplify configuration and testing choices and help structure optimized cluster configurations, Red Hat categorizes workload profiles as:

- IOPS-optimized clusters
- Throughput-optimized clusters
- Cost/capacity-optimized clusters

Table 1 provides the criteria used to identify optimal Red Hat Ceph Storage cluster configurations, their properties, and example uses. These categories are provided as general guidelines for hardware purchase and configuration decisions, and can be adjusted to satisfy unique workload blend. As the workload mix varies from organization to organization, actual hardware configurations chosen will vary.

As discussed, a single Ceph cluster can be configured to have a multiple pools to serve different workloads. For example, OSDs on IOPS-optimized servers can be configured into a pool serving MySQL workloads, while OSDs on throughput-optimized servers can be configured into a pool serving digital media performance workloads.

TABLE 1. CEPH CLUSTER OPTIMIZATION CRITERIA.

OPTIMIZATION CRITERIA	PROPERTIES	EXAMPLE USES
IOPS-OPTIMIZED	<ul style="list-style-type: none"> • Lowest cost per IOPS • Highest IOPS • Meets minimum fault domain recommendation (single server is less than or equal to 10% of the cluster) 	<ul style="list-style-type: none"> • Typically block storage • 3x replication (HDD) or 2x replication (SSD) • MySQL on OpenStack clouds
THROUGHPUT-OPTIMIZED	<ul style="list-style-type: none"> • Lowest cost per given unit of throughput • Highest throughput • Highest throughput per BTU • Highest throughput per watt • Meets minimum fault domain recommendation (single server is less than or equal to 10% of the cluster) 	<ul style="list-style-type: none"> • Block or object storage • 3x replication • Active performance storage for video, audio, and images • Streaming media
CAPACITY-OPTIMIZED	<ul style="list-style-type: none"> • Lowest cost per TB • Lowest BTU per TB • Lowest watt per TB • Meets minimum fault domain recommendation (single server is less than or equal to 15% of the cluster) 	<ul style="list-style-type: none"> • Typically object storage • Erasure coding common for maximizing usable capacity • Object archive • Video, audio, and image object archive repositories

CHOOSING A STORAGE ACCESS METHOD

Choosing a storage access method is an important design consideration. As discussed, all data in Ceph is stored in pools—regardless of data type. The data itself is stored in the form of objects via the RADOS layer (Figure 2) which:

- Avoids a single point of failure
- Provides data consistency and reliability
- Enables data replication and migration
- Offers automatic fault detection and recovery

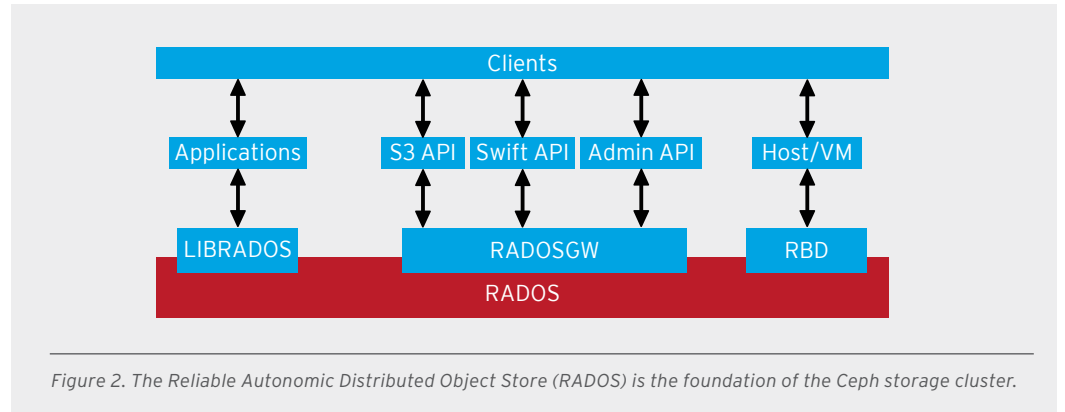


Figure 2. The Reliable Autonomic Distributed Object Store (RADOS) is the foundation of the Ceph storage cluster.

Writing and reading data in a Ceph storage cluster is accomplished using the Ceph client architecture. Ceph clients differ from competitive offerings in how they present data storage interfaces. A wide range of access methods are supported, including:

- **RADOSGW** is a bucket-based object storage gateway service with S3 compatible and OpenStack Swift compatible RESTful interfaces.
- **LIBRADOS** provides direct access to RADOS with libraries for most programming languages, including C, C++, Java™, Python, Ruby, and PHP.
- **RBD** offers a Ceph block storage device that mounts like a physical storage drive for use by both physical and virtual systems (with a Linux® kernel driver, KVM/QEMU storage backend, or user-space libraries).

Storage access method and data protection method (discussed later in this document) are interrelated. For example, Ceph block storage is currently only supported on replicated pools, while Ceph object storage is supported on either erasure-coded or replicated pools. The cost of replicated architectures is categorically more expensive than that of erasure-coded architectures due to the significant difference in media costs. Note that while CephFS distributed file storage is not yet supported on Red Hat Ceph Storage as of this writing, file systems are routinely created on top of Ceph block devices.

IDENTIFYING CAPACITY NEEDS

Identifying storage capacity may seem trivial, but it can have a distinct effect on the chosen target server architecture. In particular, storage capacity must be weighed in concert with considerations such as fault domain risk tolerance. For example, if an organization is designing a small, half-petabyte cluster, minimum server fault domain recommendations will preclude the use of ultra-dense storage servers in the architecture. Doing so avoids unacceptable fault domain risk on a small number of very large nodes. Table 2 lists broad server sizing trends, with typical types of servers categorized by both workload optimization and overall cluster size.

TABLE 2. BROAD SERVER SIZING TRENDS.

OPTIMIZATION CRITERIA	OPENSTACK STARTER (64 TB)	SMALL (250 TB)	MEDIUM (1 PB)	LARGE (2 PB)
IOPS-OPTIMIZED	<ul style="list-style-type: none"> Servers with 2-4x PCIe/NVMe slots Servers with 8-12x 2.5-inch SSD bays (SAS/SATA) 		<ul style="list-style-type: none"> NA 	<ul style="list-style-type: none"> NA
THROUGHPUT-OPTIMIZED	<ul style="list-style-type: none"> Servers with 12-16x 3.5-inch drive bays 		<ul style="list-style-type: none"> Servers with 24-36x 3.5-inch drive bays 	<ul style="list-style-type: none"> Servers with 24-36x 3.5-inch drive bays
CAPACITY-OPTIMIZED				<ul style="list-style-type: none"> Servers with 60-72x 3.5-inch drive bays

DETERMINING FAULT DOMAIN RISK TOLERANCE

It may be tempting to deploy the largest servers possible in the interest of economics. However, production environments need to provide reliability and availability for the applications they serve, and this necessity extends to the scale-out storage upon which they depend. The fault domain that a single OSD server represents is key to cluster design, so dense servers should be reserved for multi-petabyte clusters where the capacity of an individual server accounts for less than 10-15% of the total cluster capacity. This recommendation may be relaxed for less critical pilot projects. Primary factors for weighing fault domain risks include:

- **Reserving capacity for self-healing.** When a storage node fails, Ceph self-healing begins after a configured time period. The unused storage capacity of the surviving cluster nodes must be greater than the used capacity of the failed server for successful self-healing. For example, in a 10-node cluster, each node should reserve 10% unused capacity for self-healing of a failed node (in addition to reserving 10% for statistical deviation due to using algorithmic placement). As a result, each node in a cluster should operate at less than 80% of total capacity.
- **Accommodating impact on performance.** During self-healing, a percentage of cluster throughput capacity will be diverted to reconstituting object copies from the failed node on the surviving nodes. The percentage of cluster performance degradation is a function of the number of nodes in the cluster and how Ceph is configured. More nodes in the cluster results in less impact per node.

Ceph will automatically recover by re-replicating the data from the failed node using secondary copies on other nodes in the cluster. As a result, a node failure has several effects:

- Total cluster capacity is reduced by some fraction.
- Total cluster throughput is reduced by some fraction.
- The cluster enters an I/O-heavy recovery process, temporarily diverting an additional fraction of the available throughput.

The recovery process time is directly proportional to how much data was on the failed node and how much throughput the rest of the cluster can sustain. A general guide for calculating recovery time in a Ceph cluster given one disk per OSD is:

$$\text{Recovery time seconds} = (\text{disk capacity in gigabits} / \text{network speed}) / (\text{nodes} - 1)$$

For example, if a 2 TB OSD node fails in a 10-node cluster with a 10 Gb Ethernet back-end, it takes approximately three minutes for the Ceph cluster to recover with 100% of the network bandwidth and no CPU overhead. In reality, using 20% of the available 10-Gb network, it takes approximately 15 minutes to recover. With a 4 TB drive, that time will double.

Red Hat recommends the following minimum cluster sizes:

- **Supported minimum cluster size:** Three storage (OSD) servers, suitable for use cases with higher risk tolerance for performance degradation during recovery from node failure
- **Recommended minimum cluster size (IOPS- and throughput-optimized cluster):** 10 storage (OSD) servers
- **Recommended minimum cluster size (cost/capacity-optimized cluster):** 7 storage (OSD) servers

There are other considerations related to fault domain risk. Ceph replicates objects across multiple nodes in a storage cluster to provide data redundancy and higher data availability. When designing a cluster, it is important to ask these questions:

- Should the replicated node be in the same rack or multiple racks to avoid a single rack failure?
- Should the Ceph OSD traffic stay within the rack or span across racks in a dedicated or shared network?
- Are the application servers in the rack or datacenter proximate to the storage nodes?
- How many concurrent node failures can be tolerated?

Automatic, intelligent placement of object replicas across server, rack, row, and datacenter fault domains can be governed by CRUSH ruleset configuration parameters.

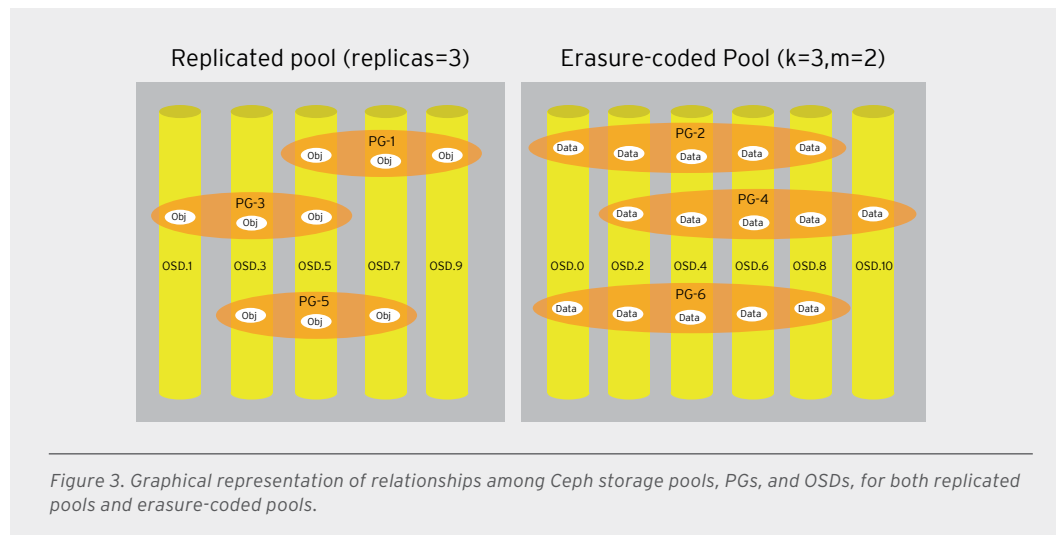
SELECTING A DATA PROTECTION METHOD

As a design decision, choosing the data protection method can affect the solution's total cost of ownership (TCO) more than any other factor. This is because the chosen data protection method strongly affects the amount of raw storage capacity that must be purchased to yield the desired amount of usable storage capacity. Applications have diverse needs for performance and availability. As a result, Ceph provides data protection at the storage pool level.

- **Replicated storage pools.** Replication makes full copies of stored objects, and is ideal for quick recovery. In a replicated storage pool, Ceph configuration defaults to a replication factor of three, involving a primary OSD and two secondary OSDs. If two of the three OSDs in a placement group become unavailable, data may be read, but write operations are suspended until at least two OSDs are operational.
- **Erasured-coded storage pools.** Erasure coding provides a single copy of data plus parity, and it is useful for archive storage and cost-effective durability and availability. With erasure coding, storage pool objects are divided into chunks using the $n=k+m$ notation, where k is the number data chunks that are created, m is the number of coding chunks that will be created to provide data protection, and n is the total number of chunks placed by CRUSH after the erasure coding process.

Ceph block storage is typically configured with 3x replicated pools and is currently not supported directly on erasure-coded pools. Ceph object storage is supported on either replicated or erasure-coded pools. Depending upon the performance needs and read/write mix of an object storage workload, an erasure-coded pool can provide an extremely cost effective solution while meeting performance requirements.

Figure 3 illustrates the relationships among Ceph storage pools, PGs, and OSDs for both replicated and erasure-coded pools. For more information on Ceph architecture, see the Ceph documentation at docs.ceph.com/docs/master/architecture/.



HARDWARE CONFIGURATION GUIDELINES

The sections that follow provide broad guidelines for the selection of both monitor nodes and OSD hosts. Actual configuration of OSD servers can vary based on application and workload optimization.

MONITOR NODES

The Ceph monitor is a datastore for the health of the entire cluster, and contains the cluster log. A minimum of three monitors are strongly recommended for a cluster quorum in production. Monitor nodes typically have fairly modest CPU and memory requirements. A 1 rack unit (1U) server with a low-cost CPU (such as an Intel Xeon processor E5-2603), 16 GB of RAM, and Gigabit Ethernet (GbE) networking should suffice in most cases. Since logs are stored on local disk(s) on the monitor node, it is important to make sure that sufficient disk space is provisioned. The monitor store should be placed on an SSD, because the leveldb store can become I/O bound.

For example, when monitoring 100 OSDs in a healthy Ceph cluster, each monitor will collect data for the 100 OSDs until all the monitors are synchronized with the same Ceph cluster information. The size of the datastore can vary, but 200 MB up to 1 or 2 GB are to be expected, depending not only on the size of the cluster, but the state change (churn) that the cluster undergoes. The logs for recovering clusters can grow quickly to over 100 GB. Abnormal monitor datastore growth should be investigated by an operator, as there is usually an underlying condition that should be remedied. Log rotation is a good practice to guarantee that available disk space is not blindly consumed. This is

particularly true if verbose debugging output is set on the monitors, since they will generate a large amount of logging information. Refer to Ceph documentation on monitor log settings for additional details. In most situations, monitors should be run on distinct nodes or on virtual machines residing on physically separate machines to prevent a single point of failure.

OSD HOSTS

Ceph OSD hosts are configured differently depending on both workload optimization and the data devices installed (HDDs, SSDs, or NVMe devices).

CPU specifications

CPU recommendations for OSD hosts differ depending on the media that is employed.

- For HDD-based OSDs, one core-GHz is recommended for each OSD. Therefore, 16 HDD-based OSDs can be supported with an 8-core 2.0GHz processor. Certain CPUs may also have custom features, such as SIMD instructions that can dramatically improve throughput in some situations.
- For NVMe-based OSDs, five core-GHz is recommended for each OSD. As such, if four OSDs are configured on a single NVMe device, a 10-core 2.0GHz processor is recommended to drive the OSDs for that single NVMe device.

Memory specifications

Red Hat typically recommends a baseline of 16 GB of RAM, with an additional 2 GB of RAM per OSD. When sizing memory requirements, it is important to consider:

- The number of OSDs per node
- The number of memory banks available
- The number of memory channels per bank
- The cost of DIMMs

Data devices

OSDs and OSD data drives are independently configurable with Ceph and OSD performance is naturally dependent on the throughput and latency of the underlying media. The actual number of OSDs configured per OSD drive depends on the type of OSD media configured on the OSD host. When using magnetic storage media, one OSD should be configured per HDD. On the other hand, an IOPS-optimized OSD host with a smaller number of high-speed SSDs might be configured with two to four OSDs per SSD to exploit the available I/O bandwidth.

Ceph uses a journal to allow it to create atomic updates, which are required to ensure data consistency. An OSD writes the data payload and metadata to a write journal before writing to the OSD's data partition. A beneficial side-effect of write journaling may be some coalescing of small writes. For performance-optimized clusters, journals are typically located on a partition of a faster media type than the OSD media. For example, a throughput-optimized OSD server typically has HDD-based OSDs, and a dedicated write journal based on an SSD.

The suggested ratio between HDD-OSDs and flash write journals are:

- Four to five HDD-OSDs for each SATA/SAS SSD journal
- 12-18 HDD-OSDs for each NVMe journal

To help decrease cost for cost/capacity-optimized clusters, journals can be co-located with OSDs on the same HDD via partitions. The OSD writes to the journal first and then copies to the OSD data partition. A write is acknowledged to the client after all OSD peers in the placement group (replicated or erasure-coded pools) have successfully written their assigned replica or shard to their write journal.

Storage Media

Performance and economics for Ceph clusters both depend heavily on an effective choice of storage media. For throughput and capacity-archive optimized clusters, magnetic media currently accounts for the bulk of the deployed storage capacity. Additionally, for throughput-optimized configurations, solid-state storage media are typically used for Ceph write journaling. For capacity-archive configurations, write journaling is co-resident on the HDDs.

- **Magnetic media.** Enterprise-, or cloud-class HDDs should be used for Ceph clusters. Desktop-class disk drives are not well suited for Ceph deployments, as they lack sufficient rotational vibration compensation for high density, high duty-cycle applications and use cases. When dozens (or hundreds) of rotating HDDs are installed in close proximity, rotational vibration quickly becomes a challenge. Failures, errors, and even overall cluster performance can be adversely affected by the rotation of neighboring disks interfering with the rapidly spinning platters in high density storage enclosures. Enterprise-class HDDs contain higher quality bearings and RV Compensation circuitry to mitigate these issues in multi-spindle applications and use cases—especially in densities above four to six HDDs in a single enclosure. Both SAS and SATA interface types are acceptable.
- **Solid state media (currently flash).** Ceph is strongly consistent storage, so every write to the Ceph cluster must be written to Ceph journals before the write is acknowledged to the client. The data remain in the journal until all replicas or shards are acknowledged as fully written. Only then will the next write happen. With SSD journals, the OSDs are able to write faster, reducing the time before a write acknowledgment is sent to the client. In some cases, several small writes can be coalesced during a single journal flush, which can also improve performance.

Key criteria to consider when selecting solid state media for Ceph include:

- **Classification.** Consumer-class SSDs should not be used; only enterprise-class SSDs should be deployed with Ceph.
- **Endurance.** Write endurance is important as Ceph write journals are heavily used and could exceed recommended program/erase cycles of an SSD rated with lower endurance. Current popular choices include devices rated at greater than 10 device writes per day (DWPD) for 5 years, translating to a lifetime total of 28 petabytes written (PBW).
- **Power fail protection.** Supercapacitors for power fail protection are vital. In the event of a power failure, supercapacitors must be properly sized to allow the drive to persist all in-flight writes to non-volatile NAND storage.
- **Performance.** For Ceph write journaling, the write throughput rating of the journal device should exceed the aggregate write throughput rating of all underlying OSD devices that are served by that journal device.

I/O controllers

Servers with either JBOD host bus adapters (HBAs) or RAID controllers can be appropriate for OSD hosts, depending on the workload and application expectations. The performance of RAID controllers varies based on the chipset and the system enclosure they are integrated into. Most vendors

have standardized on specific controller chips and either package these chips on their own branded controllers or use the chip manufacturers' controller board. For example, at the time of this writing many RAID controllers are based upon the LSI 3108 chipset.

For large block sequential I/O workload patterns, HDDs typically perform better when configured in JBOD mode. For small block random I/O patterns, configuring HDDs as single-drive RAID 0 volumes via a RAID controller with protected write-back cache can provide optimal results. Ensure that disk caches are disabled, as some RAID controllers do not do this by default.

Many modern systems that house more than eight drives have SAS expander chips on the drive hotswap backplane. Similar to network switches, SAS expanders, often allow connection of many SAS devices to a controller with a limited number of SAS lanes. Ceph node configurations with SAS expanders are well suited for large capacity-optimized clusters. However, when selecting hardware with SAS expanders, consider the impact of:

- Adding extra latency
- Oversubscribed SAS lanes
- STP overhead of tunneling SATA over SAS

Because of backplane oversubscription or poor design, some servers used for Ceph deployments have encountered sub-par performance in systems that use SAS expanders, although it is not universally the case. The type of controller, expander, and even brand of drive and firmware all play a part in determining performance.

Network interfaces

Providing sufficient network bandwidth is essential for an effective and performant Ceph cluster. Fortunately, network technology is improving rapidly with standard Ethernet-based interfaces available with ever-increasing bandwidth. In servers employed as OSD hosts, network capacity should generally relate to storage capacity. For smaller OSD hosts with 12-16 drive bays, 10 Gigabit Ethernet is typically sufficient. For larger OSD hosts with 24-72 drive bays, 40 Gigabit Ethernet may be preferred to provide the required bandwidth and throughput.

Physical deployment characteristics must also be taken into account. If the nodes are spread across multiple racks in the datacenter, then the network design should ensure high bisectional bandwidth, and minimal network diameter. Ideally, each OSD server should have two network interfaces for data traffic—one connected to the client systems, and one for the private network connecting the OSD servers.

BROAD OSD HOST CONFIGURATION TRENDS

Taking all of these considerations into account, Table 3 provides general guidance for configuring OSD hosts for Red Hat Ceph Storage.

TABLE 3. BROAD SERVER CONFIGURATION TRENDS.

OPTIMIZATION CRITERIA	OPENSTACK STARTER (100 TB)	SMALL (250 TB)	MEDIUM (1 PB)	LARGE (2 PB)
IOPS-OPTIMIZED	<ul style="list-style-type: none"> • Ceph RBD (block) pools • Either: OSDs on four SAS/SATA SSDs with one NVMe SSD write journal; or OSDs on 1-4 NVMe SSDs, with journals co-located on each device* • 20 CPU core-GHz per NVMe SSD with four OSDs per device • 6 CPU core-GHz per SATA/SAS SSD with two OSDs per device • Data protection: Replication (2x on SSD-based OSDs) with regular backups to the object storage pool • 2-4 OSDs per SSD or NVMe drive 		• NA	• NA
THROUGHPUT-OPTIMIZED	<ul style="list-style-type: none"> • Ceph RBD (block) or Ceph RGW (object) pools • OSDs on HDDs with dedicated write journals (4-5:1 ratio of HDDs with SSDs, or 12-18:1 ratio of HDDs with NVMe) • One CPU core-GHz per OSD • Mid-bin, dual-socket CPU (single-socket adequate for servers <= 12 OSDs) (e.g. Intel Xeon E5-2630 v3 for 12 HDD OSDs, and dual Intel Xeon E5-2650 v3 for 36 HDD OSDs) • Data protection: Replication (read-intensive or mixed read/write) or erasure-coded (write-intensive) • High-bandwidth networking, greater than 10 GbE for servers with more than 12-16 drives 			
CAPACITY-OPTIMIZED	• NA	<ul style="list-style-type: none"> • Ceph RGW (object) pools • OSD HDDs with Ceph write journals co-located on HDDs • Mid-bin, dual-socket CPU (single-socket adequate for servers <= 12 OSDs) • Data protection: Erasure-coded 		

* All SSDs should be enterprise-class, meeting the requirements noted above.

CONCLUSION

Selecting the right hardware for target workloads can be a challenge, and this is especially true for software-defined storage solutions like Ceph, that run on commodity hardware. Because every environment differs, the general guidelines for sizing CPU, memory, and disk per node in this document should be mapped to a preferred vendor's product portfolio for determining appropriate server hardware. Additionally, the guidelines and best practices highlighted in this document are not a substitute for running baseline benchmarks before going into production.

Red Hat has conducted extensive testing with a number of vendors who supply hardware optimized for Ceph workloads. For specific information on selecting servers for running Red Hat Ceph Storage, please refer to the tested configurations documented in the "Red Hat Ceph Storage Hardware Selection Guide." Detailed information and Red Hat Ceph Storage test results can be found in performance and sizing guides for popular hardware vendors.



ABOUT RED HAT

Red Hat is the world's leading provider of open source solutions, using a community-powered approach to provide reliable and high-performing cloud, virtualization, storage, Linux, and middleware technologies. Red Hat also offers award-winning support, training, and consulting services. Red Hat is an S&P company with more than 80 offices spanning the globe, empowering its customers' businesses.



facebook.com/redhatinc
@redhatnews
linkedin.com/company/red-hat

redhat.com
#INC0387897_0416

NORTH AMERICA
1 888 REDHAT1

**EUROPE, MIDDLE EAST,
AND AFRICA**
00800 7334 2835
europe@redhat.com

ASIA PACIFIC
+65 6490 4200
apac@redhat.com

LATIN AMERICA
+54 11 4329 7300
info-latam@redhat.com

Copyright © 2016 Red Hat, Inc. Red Hat, Red Hat Enterprise Linux, the Shadowman logo, and JBoss are trademarks of Red Hat, Inc., registered in the U.S. and other countries. The OpenStack® Word Mark and OpenStack Logo are either registered trademarks / service marks or trademarks / service marks of the OpenStack Foundation, in the United States and other countries and are used with the OpenStack Foundation's permission. We are not affiliated with, endorsed or sponsored by the OpenStack Foundation or the OpenStack community. Linux® is the registered trademark of Linus Torvalds in the U.S. and other countries.