



# NVIDIA VIDEO DECODER (NVDEC) INTERFACE

PG-08085-001\_v07 | June 2016

**Programming Guide**



## DOCUMENT CHANGE HISTORY

PG-08085-001\_v07

Version	Date	Authors	Description of Change
1.0	2016/6/10	VU/CC	Initial release

# TABLE OF CONTENTS

<b>Chapter 1. Overview</b> .....	1
1.1 Decoding & Interoperability .....	1
1.2 Support .....	2
<b>Chapter 2. Video Decoder Capabilities</b> .....	1
<b>Chapter 3. Video Playback and Decoder Pipeline</b> .....	3
3.1 Decoder Pipeline .....	4
<b>Chapter 4. NVIDIA Video Decoder (NVDEC) API</b> .....	6
4.1 Video Decoder APIs .....	6
4.2 Creating a Decoder .....	7
4.3 Decoding Surfaces .....	9
4.4 Processing and Displaying Frames .....	11
4.5 Writing an Efficient Decode-Display Application .....	13

## LIST OF FIGURES

Figure 1. Video decoder pipeline using NVDECODE API .....	4
---	---

## LIST OF TABLES

Table 1. Hardware Video Decoder Capabilities.....	1
---	---

# Chapter 1.

## OVERVIEW

The NVIDIA Video Decoder Interface hereafter referred to as NVDECODE APIs lets developers access the video decoding features of NVIDIA graphics hardware and also interoperates video with compute and graphics.

### 1.1 DECODING & INTEROPERABILITY

Compressed video streams are decoded directly to video memory. With frames in video memory, post processing can be done using CUDA. Additionally clients can use NVIDIA CUDA APIs for synchronous/asynchronous memory transfers between video memory and system memory. Decoded video frames can either be presented to the display with graphics interoperability for video playback, or frames can be passed directly to a dedicated hardware encoder (NVENC) for video transcoding.

## 1.2 SUPPORT

The API is supported on multiple OS platforms<sup>1</sup> and works in conjunction with NVIDIA's CUDA, graphics, and encoder capabilities. The NVDECODE API supports the following video codec formats:

- ▶ MPEG-2,
- ▶ VC-1,
- ▶ H.264 (AVCHD),
- ▶ H.265 (HEVC),
- ▶ VP8 and
- ▶ VP9 (profile 0).

Refer to Chapter 2 for complete details about the video capabilities for each GPU architecture.

---

<sup>1</sup> Windows and Linux OS

# Chapter 2. VIDEO DECODER CAPABILITIES

Table 1 shows the codec support and capabilities of the hardware video decoder for each GPU architecture.

Table 1. Hardware Video Decoder Capabilities

GPU Architecture	MPEG-2	VC-1	H.264/AVCHD	H.265/HEVC	VP8	VP9
Fermi (GF1xx)	<u>Maximum Resolution:</u> 4080x4080	<u>Maximum Resolution:</u> 2048x1024 1024x2048	<u>Maximum Resolution:</u> 4096x4096 <u>Profile:</u> Baseline, Main, High profile up to Level 4.1	Unsupported	Unsupported	Unsupported
Kepler (GK1xx)	<u>Maximum Resolution:</u> 4080x4080	<u>Maximum Resolution:</u> 2048x1024 1024x2048	<u>Maximum Resolution:</u> 4096x4096 <u>Profile:</u> Main, High profile up to Level 4.1	Unsupported	Unsupported	Unsupported
Maxwell Gen 1 (GM10x)	<u>Maximum Resolution:</u> 4080x4080	<u>Maximum Resolution:</u> 2048x1024 1024x2048	<u>Maximum Resolution:</u> 4096x4096 <u>Profile:</u> Baseline, Main, High profile up to Level 5.1	Unsupported	Unsupported	Unsupported

GPU Architecture	MPEG-2	VC-1	H.264/AVCHD	H.265/HEVC	VP8	VP9
Maxwell Gen 2 (GM20x)	<u>Maximum Resolution:</u> 4080x4080	<u>Maximum Resolution:</u> 2048x1024 1024x2048  Max bitrate: 60Mbps	<u>Maximum Resolution:</u> 4096x4096 <u>Profile:</u> Baseline, Main, High profile up to Level 5.1	Unsupported	<u>Maximum Resolution:</u> 4096x4096	Unsupported
Maxwell Gen 2 (GM206)	<u>Maximum Resolution:</u> 4080x4080	<u>Maximum Resolution:</u> 2048x1024 1024x2048  Interlaced	<u>Maximum Resolution:</u> 4096x4096 <u>Profile:</u> Baseline, Main, High profile up to Level 5.1	<u>Maximum Resolution:</u> 4096x2304 <u>Profile:</u> Main profile up to Level 5.1	<u>Maximum Resolution:</u> 4096x4096	<u>Maximum Resolution:</u> 4096x2304 <u>Profile:</u> Profile 0
GP100	<u>Maximum Resolution:</u> 4080x4080	<u>Maximum Resolution:</u> 2048x1024 1024x2048	<u>Maximum Resolution:</u> 4096x4096 <u>Profile:</u> Baseline, Main, High profile up to Level 5.1	<u>Maximum Resolution:</u> 4096x4096 <u>Profile:</u> Main profile up to Level 5.1	<u>Maximum Resolution:</u> 4096x4096	<u>Maximum Resolution:</u> <u>4096x4096</u> <u>Profile:</u> <u>Profile 0</u>
GP10x	<u>Maximum Resolution:</u> 4080x4080	<u>Maximum Resolution:</u> 2048x1024 1024x2048	<u>Maximum Resolution:</u> 4096x4096 <u>Profile:</u> Baseline, Main, High profile up to Level 5.1	<u>Maximum Resolution:</u> 8192x8192 <u>Profile:</u> Main profile up to Level 5.1	Supported <sup>2</sup> <u>Maximum Resolution:</u> 4096x4096	<u>Maximum Resolution:</u> 8192x8192 <u>Profile:</u> Profile 0

---

<sup>2</sup> Supported only on GP104

# Chapter 3. VIDEO PLAYBACK AND DECODER PIPELINE

Sample<sup>3</sup> applications `NvDecodeD3D9` (DirectX 9), `NvDecodeD3D11` (DirectX 11) `NvDecodeGL` (OpenGL on Windows and Linux), included in the SDK package, demonstrate the following functions in video playback:

1. Parse the video input source.
2. Decode video on GPU using NVDECODE API.
3. Convert decoded surface NV12 format to RGBA.
4. Map RGBA surface to DirectX 9.0 or OpenGL surface.
5. Draw texture to screen.

Sample application `NvTranscoder` included in the SDK package demonstrates how to set up an end-to-end vide transcode pipeline using NVDECODE and NVENCODE APIs, with following functions:

1. Parse the video input source.
2. Decode video on GPU using NVDECODE API.
3. Send YUV video frame to encoding using the NVENCODE API.
4. Receive a compressed video bitstream back to the host.

---

<sup>3</sup> Location: `./Samples` in the Video Codec SDK package

## 3.1 DECODER PIPELINE

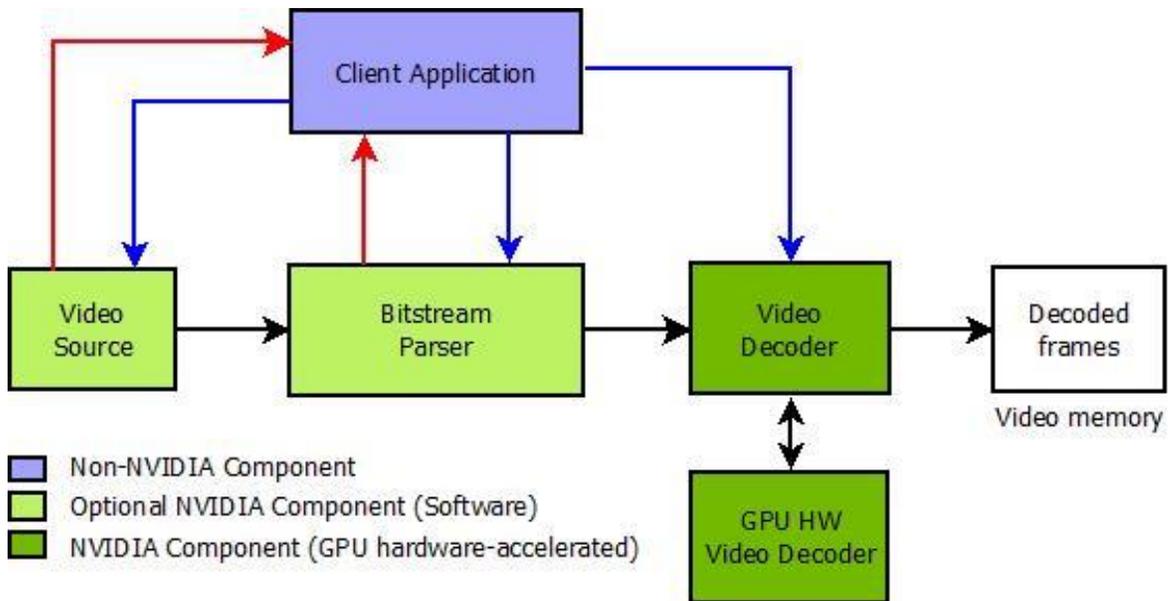


Figure 1. Video decoder pipeline using NVDECODE API

Figure 1 shows the decoder pipeline using NVDECODE API. The solid black lines show the data flow between modules. The solid colored lines represent process flow.

1. The application thread (referred to as the primary thread) calls `cvidCreateVideoSource()`, which spawns a de-multiplexer thread (referred to as the secondary thread).
2. The primary thread (colored in blue above) calls `cvidCreateVideoParser()` to create the parser. It also creates the decoder by calling `cvidCreateDecoder()`.
3. The secondary thread (colored in red above) makes the following callbacks given that the function pointers are not NULL. The callbacks are serial:
  - a) Handle video data: The callback implementation calls `cvidParseVideoData()` to parse the video data.
  - b) Handle video sequence: The callback is made when there is sequence change.
  - c) Handle picture decode: The callback implementation calls `cvidDecodePicture()` to decode the frame.
  - d) Handle Picture Display: The callback implementation signals the primary thread to display the picture.

4. The primary thread calls `cuidMapVideoFrame()` to get the pitch and CUDA device pointer to the surface which contains the decoded/post-processed frame. Thereafter it calls `cuidUnmapVideoFrame()` as the complimentary operation.
5. The primary thread destroys the resources by calling `cuidDestroyDecoder()`, `cuidDestroyVideoParser()` and `cuidDestroyVideoSource()`.

The sample applications use all three components - Video Source, Video Parser, and Video Decoder. The components are not dependent on each other and hence can be used independently. The user can unplug the Video Source and Video Parser and plug-in an implementation of his own. In this document we will be concentrating particularly on the Video Decoder (colored dark green in Figure 1) and the stages following decode (format conversion and display using OpenGL or DirectX). It is highly recommended that the user use his own implementation for Video Source and Video Parser. These two components are not hardware-accelerated or optimized and users may want to have their own customized parsers which have better performance.

# Chapter 4. NVIDIA VIDEO DECODER (NVDECODE) API

The NVDECODE API consists of two main header-files: `dynlink_cuviddec.h` and `dynlink_nvcuvid.h`. The samples in NVIDIA Video Codec SDK dynamically load the library functions and only include `dynlink_cuviddec.h` and `dynlink_nvcuvid.h` in the source files. These headers can be found under `./Samples/common/inc` folder in the Video Codec SDK package. The Windows DLL `nvcuvid.dll` is included in the NVIDIA display driver for Windows. The Linux library `libnvcuvid.so` is included with NVIDIA display driver for Linux.

## 4.1 VIDEO DECODER APIS

The Video Decoder consists of the following APIs:

```
// Create the Decoder Object
CUresult cuvidCreateDecoder(CUvideodecoder *phDecoder,
                           CUIDEDECODERCREATEINFO *pdci);
// Destroy the Decoder Object
CUresult cuvidDestroyDecoder(CUvideodecoder hDecoder);

// Decode a single picture (field or frame)
CUresult cuvidDecodePicture(CUvideodecoder hDecoder,
                           CUIDEDECODERCREATEINFO *pPicParams);

// Post-Process and map a video frame for use in CUDA
CUresult cuvidMapVideoFrame(CUvideodecoder hDecoder, int PicIdx,
                           unsigned int* pDevPtr, unsigned int*
                           pPitch, CUIDEDECODERCREATEINFO* pVPP);
```

```
// Unmap the previously mapped video frame
CUresult cuvidUnmapVideoFrame(CUvideodecoder hDecoder, unsigned int
DevPtr
```

## 4.2 CREATING A DECODER

The sample application uses the API `cuvidCreateDecoder()` through a C++ wrapper class `VideoDecoder` defined in `VideoDecoder.h`. The class's constructor is a good starting point to see how to set up the `CUVIDDECODECREATEINFO` for the `cuvidCreateDecoder()` method. Most importantly, the structure `CUVIDDECODECREATEINFO` contains the following information about the stream to be decoded:

- ▶ Codec type
- ▶ Frame size
- ▶ Chroma format

The user also specifies various properties of the output that the decoder should generate:

- ▶ Output surface format (currently only NV12 is supported)
- ▶ Output frame size
- ▶ Maximum number of output surfaces: This is the maximum number of surfaces that the client code will simultaneously map for display.

The user also needs to specify the maximum number of surfaces the decoder may allocate for decoding.

The following pseudo-code demonstrates the setup of decoder in case of scaling, cropping, or aspect ratio conversion.

```

// Scaling. Source size is 1280x960. Scale to 1920x1080.
CUresult rResult;
unsigned int uScaleW, uScaleH;
uScaleW = 1920;
uScaleH = 1080;
...

CUVIDDECODECREATEINFO stDecodeCreateInfo;
memset(&stDecodeCreateInfo, 0, sizeof(CUVIDDECODECREATEINFO));

... // setup the structure members

stDecodeCreateInfo.ulTargetWidth = uScaleWidth;
stDecodeCreateInfo.ulTargetHeight = uScaleHeight;

rResult = cuvidCreateDecoder(&hDecoder, &stDecodeCreateInfo);
...

```

```

// Cropping. Source size is 1280x960
CUresult rResult;
unsigned int uCropL, uCropR, uCropT, uCropB;
uCropL = 30;
uCropR = 700;
uCropT = 20;
uCropB = 500;
...

CUVIDDECODECREATEINFO stDecodeCreateInfo;
memset(&stDecodeCreateInfo, 0, sizeof(CUVIDDECODECREATEINFO));

... // setup structure members

stDecodeCreateInfo.display_area.left = uCropL;
stDecodeCreateInfo.display_area.right = uCropR;
stDecodeCreateInfo.display_area.top = uCropT;
stDecodeCreateInfo.display_are.bottom = uCropB;

rResult = cuvidCreateDecoder(&hDecoder, &stDecodeCreateInfo);
...

```

```

// Aspect Ratio Conversion. Source size is 1280x960(4:3). Convert to
// 16:9
CUresult rResult;
unsigned int uCropL, uCropR, uCropT, uCropB;
uDispAR_L = 0;
uDispAR_R = 1280;
uDispAR_T = 70;
uDispAR_B = 790;
...

CUVIDDECODECREATEINFO stDecodeCreateInfo;
memset(&stDecodeCreateInfo, 0, sizeof(CUVIDDECODECREATEINFO));

... // setup structure members

stDecodeCreateInfo.target_rect.left    = uDispAR_L;
stDecodeCreateInfo.target_rect.right   = uDispAR_R;
stDecodeCreateInfo.target_rect.top     = uDispAR_T;
stDecodeCreateInfo.target_rect.bottom  = uDispAR_B;

reResult = cuvidCreateDecoder(&hDecoder, &stDecodeCreateInfo);
...

```

## 4.3 DECODING SURFACES

The classes `VideoSource` and `VideoParser` wrap the calls to Video Source and Video Parser components of Figure 1. The `VideoParser` class implements three callback functions, two of which are explained below:

```

// called by the video parser to decode a single picture. Since the
// parser will deliver data as fast as it can, we need to make sure
// that the picture index we're attempting to use for decode is no
// longer used for display.
static int CUDAAPI HandlePictureDecode(void *pUserData,
                                       CUVIDPICPARAMS *pPicParams);

// called by the video parser to display a video frame (in case of
// field pictures, there may be two decode calls per one display call,
// since two fields make up one frame)
static int CUDAAPI HandlePictureDisplay(void *pUserData,
                                       CUVIDPARSERDISPINFO *pPicParams);

```

`VideoParser` passes a `CUVIDPICPARAMS` structure to the callback which can be passed without any modifications to the function `cuvidDecodePicture()`. The `CUVIDPICPARAMS` structure contains all the information necessary for the decoder to decode a frame or field. In particular, it contains pointers to the video bitstream,

information about frame size, flags denoting whether it's a field or a frame, bottom or top field, etc.

The decoded result gets associated with a picture-index value in the `CUVIDPICPARAMS` structure, which is also provided by the parser. This picture index is later used to map the decoded frames to CUDA memory.

The implementation of `HandlePictureDecode()` in the sample application waits if the output queue is full. When a slot in the queue becomes available, it simply invokes the `cuvidDecodePicture()` function, passing the `pPicParams` as received from the parser.

The `HandlePictureDisplay()` method is passed a `CUVIDPARSERDISPINFO` structure which contains the necessary data for displaying a frame; i.e. frame index of the decoded frame (as given to the decoder), and some information relevant for display such as frame time, field information, etc. The parser calls this method for frames in the order that they should be displayed.

The implementation of `HandlePictureDisplay()` method in the sample application simply enqueues the `pPicParams` passed by the parser into the `FrameQueue` object.

The `FrameQueue` is used to implement a producer-consumer pattern for passing frames (or better, references to decoded frames) between the `VideoSource`'s decoding thread and the application's main thread, which is responsible for displaying them on the screen.

## 4.4 PROCESSING AND DISPLAYING FRAMES

The user needs to call `cuidmapVideoFrame()` to get the CUDA device pointer and pitch of the surface that has the decoded frame. The following is a pseudo-code that demonstrates using `cuidmapVideoFrame()` and `cuidunmapVideoFrame()`.

```
// MapFrame: Call cuidmapVideoFrame and get the devptr and associated
// pitch. Copy this surface (in device memory) to host memory using
// CUDA device to host memcpy.

bool MapFrame()
{
    CUVIDPARSEDISPINFO stDispInfo;
    CUVIDPROCPARAMS stProcParams;
    CUresult rResult;
    unsigned int cuDevPtr; int nPitch, nPicIdx;
    unsigned char* pHostPtr;

    memset(&stDispInfo, 0, sizeof(CUVIDPARSEDISPINFO));
    memset(&stProcParams, 0, sizeof(CUVIDPROCPARAMS));

    ... // setup stProcParams if required

    // retrieve the frames from the Frame Display Queue. This Queue is
    // is populated in HandlePictureDisplay.
    if (g_pFrameQueue->dequeue(&stDispInfo))
    {
        nPicIdx = stDispInfo.picture_index;
        rResult = cuidmapVideoFrame(&hDecoder, nPicIdx, &cuDevPtr,
                                   &nPitch, &stProcParams);

        // use CUDA based Device to Host memcpy
        pHostPtr = cuMemAllocHost((void** )&pHostPtr, nPitch);
        if (pHostPtr)
        {
            rResult = cuMemcpyDtoH(pHostPtr, cuDevPtr, nPitch);
        }
        rResult = cuidunmapVideoFrame(&hDecoder, cuDevPtr);
    }

    ... // Dump YUV to a file

    if (pHostPtr)
    {
        cuMemFreeHost(pHostPtr);
    }
    ...
}
```

The function `copyDecodedFrameToTexture()` in `videoDecode.cpp` does something more than the above pseudo-code. It retrieves the frame (decoded surface) from `FrameQueue` as above. Uses `cuidMapVideoFrame()` to get the CUDA device pointer and the associated pitch of the decoded surface. It maps a D3D/OGL texture to be used by CUDA (interop surface). It then calls `cudaPostProcessFrame()` to do the color space conversion from NV12 to RGBA. The texture holds the RGBA surface. This texture can now be drawn to the screen.

The following list summarizes the function calls involved (refer sample apps) in the display and post-process pipeline:

1. `cuidMapVideoFrame` – gets a CUDA device pointer from decoded frame of a Video Decoder (using map).
2. `cuD3D9ResourceGetMappedPointer` – For `cudaDecodeD3D9`, this function retrieves a CUDA device pointer from a D3D9 texture.
3. `cuGLMapBufferObject` – For `cudaDecodeGL`, this function retrieves a CUDA device pointer from an OpenGL PBO (Pixel Buffer Object).
4. `cudaPostProcessFrame` – calls all subsequent CUDA post-process functions on that frame, and writes the result directly to the Mapped D3D texture.
5. `cuD3D9UnmapResources` – For `NvDecodeD3D9`, the CUDA driver will release the pointer back to the D3D9 driver. This tells the Direct3D driver that CUDA is finished modifying the resource, and that it is safe to use it with D3D9.
6. `cuGLUnmapBufferObject` – For `NvDecodeGL`, the CUDA driver will release the pointer back to the OpenGL driver. This tells the OpenGL driver that CUDA is finished modifying the resource, and that it is safe to use it with OpenGL.
7. `cuidUnmapVideoFrame` – Unmap the previously mapped frame.

## 4.5 WRITING AN EFFICIENT DECODE-DISPLAY APPLICATION

The NVDEC engine on NVIDIA GPUs is a dedicated hardware block, which decodes the input video bitstream in supported formats. A typical video decode and display application consists of the following stages:

1. Video bitstream parser
2. Video decoder
3. Post-processor
4. Screen display

Of these, post-processing (such as scaling, color space conversion, noise reduction, color enhancement etc.) can be effectively performed using user-defined CUDA kernels.

The post-processed frames can then be sent to the display engine for displaying on the screen, if required. Note that this operation is outside the scope of NVDECODE APIs.

The sample applications included with the Video Codec SDK are written to demonstrate the functionality of various APIs but they are by no means fully optimized applications. In fact, programmers are strongly encouraged to ensure that their application is well-designed, with various stages in the decode-postprocess-display pipeline structured in an efficient manner to achieve desired performance.

As a starting point, an optimized implementation may make use of independent threads for bitstream decode and display as follows:

1. **Decode Thread:** This thread calls `cuvvidDecodePicture()` and pushes the decoded frame to the display queue. This continues as long as there are frames to decode.
2. **Display thread:** This thread reads the display queue and checks if there are any decoded frames. If yes, then it calls `cuvvidMapVideoFrame()` to get the CUDA device pointer and pitch of the frame. The resulting CUDA device pointer can be used for CUDA post-processing of the decoded video frames using user-defined CUDA kernels. Finally, it is necessary to call `cuvvidUnMapVideoFrame()` so that the decoded frame buffer is unmapped by the driver. This continues as long as there are decoded frames in the display queue and end of decode has not been reached. The display thread presents the contents of the post processed video frame to an OpenGL or Direct3D surface, using CUDA interoperability.

To ensure that the video frames will playback without stuttering or hitching, it is necessary to ensure that decode and display threads do not get blocked. Two or more D3D9/D3D11 or OpenGL surfaces allows double or triple buffered playback. This allows both decode and display to run on different surfaces without being blocked.

## **Notice**

ALL NVIDIA DESIGN SPECIFICATIONS, REFERENCE BOARDS, FILES, DRAWINGS, DIAGNOSTICS, LISTS, AND OTHER DOCUMENTS (TOGETHER AND SEPARATELY, "MATERIALS") ARE BEING PROVIDED "AS IS." NVIDIA MAKES NO WARRANTIES, EXPRESSED, IMPLIED, STATUTORY, OR OTHERWISE WITH RESPECT TO THE MATERIALS, AND EXPRESSLY DISCLAIMS ALL IMPLIED WARRANTIES OF NONINFRINGEMENT, MERCHANTABILITY, AND FITNESS FOR A PARTICULAR PURPOSE.

Information furnished is believed to be accurate and reliable. However, NVIDIA Corporation assumes no responsibility for the consequences of use of such information or for any infringement of patents or other rights of third parties that may result from its use. No license is granted by implication of otherwise under any patent rights of NVIDIA Corporation. Specifications mentioned in this publication are subject to change without notice. This publication supersedes and replaces all other information previously supplied. NVIDIA Corporation products are not authorized as critical components in life support devices or systems without express written approval of NVIDIA Corporation.

## **HDMI**

HDMI, the HDMI logo, and High-Definition Multimedia Interface are trademarks or registered trademarks of HDMI Licensing LLC.

## **OpenCL**

OpenCL is a trademark of Apple Inc. used under license to the Khronos Group Inc.

## **Trademarks**

NVIDIA and the NVIDIA logo are trademarks or registered trademarks of NVIDIA Corporation in the U.S. and other countries. Other company and product names may be trademarks of the respective companies with which they are associated.

## **Copyright**

© 2016 NVIDIA Corporation. All rights reserved.