

# NVIDIA CUDA TOOLKIT 10.1.243

RN-06722-001 \_v10.1 | August 2019

Release Notes for Windows, Linux, and Mac OS

### TABLE OF CONTENTS

Chapter 1. CUDA Toolkit Major Components1
Chapter 2. CUDA 10.1 Update 2 Release Notes
2.1. General CUDA
2.1.1. CUDA Tools
2.1.1.1. CUDA Compilers
2.2. CUDA Libraries
2.2.1. cuBLAS Library
2.2.2. cuSOLVER Library
2.2.3. cuSPARSE Library
2.2.4. NPP Library5
2.2.5. nvJPEG Library5
2.3. Deprecated Features5
2.4. Resolved Issues
2.4.1. CUDA Tools
2.4.2. CUDA Compilers
2.4.3. CUDA Libraries
2.5. Known Issues7
2.5.1. General CUDA7
2.5.2. CUDA Tools
Chapter 3. CUDA 10.1 Update 1 Release Notes
3.1. General CUDA
3.2. CUDA Tools
3.2.1. CUDA Compilers
3.3. CUDA Libraries
3.3.1. cuBLAS Library10
3.3.2. cuSPARSE Library10
3.3.3. cuFFT Library
3.3.4. NPP Library 10
3.3.5. nvJPEG Library 11
3.4. Deprecated Features11
3.5. Resolved Issues 11
3.5.1. General CUDA11
3.5.2. CUDA Compilers 11
3.5.3. CUDA Profiler 12
3.5.4. CUDA Libraries12
3.5.5. CUDA Tools
3.6. Known Issues
3.6.1. General CUDA13
3.6.2. CUDA Tools
3.6.3. CUDA Libraries14

Chapter 4. CUDA 10.1 Release Notes
4.1. General CUDA15
4.2. CUDA Tools
4.2.1. CUDA Compilers
4.2.2. CUDA Profiler
4.2.3. CUDA-MEMCHECK17
4.3. CUDA Libraries
4.3.1. cuBLAS Library17
4.3.2. cuSOLVER Library18
4.3.3. cuSPARSE Library
4.3.4. cuFFT Library
4.3.5. cuRAND Library
4.3.6. NPP Library 19
4.3.7. nvJPEG Library 19
4.4. Deprecated Features19
4.5. Resolved Issues
4.5.1. CUDA Compilers
4.5.2. CUDA Libraries20
4.6. Known Issues
4.6.1. General CUDA21
4.6.2. CUDA Tools
4.6.3. CUDA Libraries22
Chapter 5. Thrust v1.9.4 Release Notes 23
5.1. New Features
5.2. New Examples
5.3. Other Enhancements27
5.3.1. Tagged Pointer Enhancements
5.3.2. Iterator Enhancements27
5.3.3. Testing Enhancements27
5.4. Resolved Issues
Chapter 6. CUDA Tegra Release Notes
6.1. New Features
6.2. Known Issues and Limitations
6.3. Resolved Issues
6.4. Deprecated Issues

### LIST OF TABLES

Table 1 CUDA Toolkit and Compatible Driver versions
---

# Chapter 1. CUDA TOOLKIT MAJOR COMPONENTS

This section provides an overview of the major components of the CUDA Toolkit and points to their locations after installation.

#### Compiler

The CUDA-C and CUDA-C++ compiler, **nvcc**, is found in the **bin**/ directory. It is built on top of the NVVM optimizer, which is itself built on top of the LLVM compiler infrastructure. Developers who want to target NVVM directly can do so using the Compiler SDK, which is available in the **nvvm**/ directory.

Please note that the following files are compiler-internal and subject to change without any prior notice.

- any file in include/crt and bin/crt
- include/common\_functions.h, include/device\_double\_functions.h, include/device\_functions.h, include/host\_config.h, include/ host defines.h, and include/math functions.h
- nvvm/bin/cicc
- bin/cudafe++, bin/bin2c, and bin/fatbinary

#### Tools

The following development tools are available in the **bin/** directory (except for Nsight Visual Studio Edition (VSE) which is installed as a plug-in to Microsoft Visual Studio, Nsight Compute and Nsight Systems are available in a separate directory).

- IDEs: nsight (Linux, Mac), Nsight VSE (Windows)
- Debuggers: cuda-memcheck, cuda-gdb (Linux), Nsight VSE (Windows)
- Profilers: Nsight Systems, Nsight Compute, nvprof, nvvp, Nsight VSE (Windows)
- Utilities: cuobjdump, nvdisasm, gpu-library-advisor

#### Libraries

The scientific and utility libraries listed below are available in the **lib**/ directory (DLLs on Windows are in **bin**/), and their interfaces are available in the **include**/ directory.

- cublas (BLAS)
- cublas\_device (BLAS Kernel Interface)

- cuda\_occupancy (Kernel Occupancy Calculation [header file implementation])
- cudadevrt (CUDA Device Runtime)
- cudart (CUDA Runtime)
- cufft (Fast Fourier Transform [FFT])
- cupti (CUDA Profiling Tools Interface)
- **curand** (Random Number Generation)
- **cusolver** (Dense and Sparse Direct Linear Solvers and Eigen Solvers)
- cusparse (Sparse Matrix)
- nvJPEG (JPEG encoding/decoding)
- npp (NVIDIA Performance Primitives [image and signal processing])
- nvblas ("Drop-in" BLAS)
- nvcuvid (CUDA Video Decoder [Windows, Linux])
- nvgraph (CUDA nvGRAPH [accelerated graph analytics])
- nvml (NVIDIA Management Library)
- nvrtc (CUDA Runtime Compilation)
- nvtx (NVIDIA Tools Extension)
- thrust (Parallel Algorithm Library [header file implementation])

#### **CUDA Samples**

Code samples that illustrate how to use various CUDA and library APIs are available in the samples/ directory on Linux and Mac, and are installed to C:\ProgramData \NVIDIA Corporation\CUDA Samples on Windows. On Linux and Mac, the samples/ directory is read-only and the samples must be copied to another location if they are to be modified. Further instructions can be found in the *Getting Started Guides* for Linux and Mac.

#### Documentation

The most current version of these release notes can be found online at http:// docs.nvidia.com/cuda/cuda-toolkit-release-notes/index.html. Also, the **version.txt** file in the root directory of the toolkit will contain the version and build number of the installed toolkit.

Documentation can be found in PDF form in the doc/pdf/ directory, or in HTML form at doc/html/index.html and online at http://docs.nvidia.com/cuda/ index.html.

#### **CUDA Driver**

Running a CUDA application requires the system with at least one CUDA capable GPU and a driver that is compatible with the CUDA Toolkit. See Table 1. For more information various GPU products that are CUDA capable, visit https://developer.nvidia.com/cuda-gpus.

Each release of the CUDA Toolkit requires a minimum version of the CUDA driver. The CUDA driver is backward compatible, meaning that applications compiled against a particular version of the CUDA will continue to work on subsequent (later) driver releases.

More information on compatibility can be found at https://docs.nvidia.com/cuda/ cuda-c-best-practices-guide/index.html#cuda-runtime-and-driver-api-version.

CUDA Toolkit	Linux x86_64 Driver Version	Windows x86_64 Driver Version		
CUDA 10.1 (10.1.105 general release, and updates)	>= 418.39	>= 418.96		
CUDA 10.0.130	>= 410.48	>= 411.31		
CUDA 9.2 (9.2.148 Update 1)	>= 396.37	>= 398.26		
CUDA 9.2 (9.2.88)	>= 396.26	>= 397.44		
CUDA 9.1 (9.1.85)	>= 390.46	>= 391.29		
CUDA 9.0 (9.0.76)	>= 384.81	>= 385.54		
CUDA 8.0 (8.0.61 GA2)	>= 375.26	>= 376.51		
CUDA 8.0 (8.0.44)	>= 367.48	>= 369.30		
CUDA 7.5 (7.5.16)	>= 352.31	>= 353.66		
CUDA 7.0 (7.0.28)	>= 346.46	>= 347.62		

Table 1	CUDA	Toolkit	and	Com	patible	Driver	Versions

For convenience, the NVIDIA driver is installed as part of the CUDA Toolkit installation. Note that this driver is for development purposes and is not recommended for use in production with Tesla GPUs.

For running CUDA applications in production with Tesla GPUs, it is recommended to download the latest driver for Tesla GPUs from the NVIDIA driver downloads site at http://www.nvidia.com/drivers.

During the installation of the CUDA Toolkit, the installation of the NVIDIA driver may be skipped on Windows (when using the interactive or silent installation) or on Linux (by using meta packages).

For more information on customizing the install process on Windows, see http://docs.nvidia.com/cuda/cuda-installation-guide-microsoft-windows/index.html#install-cuda-software.

For meta packages on Linux, see https://docs.nvidia.com/cuda/cuda-installation-guide-linux/index.html#package-manager-metas

#### **CUDA-GDB Sources**

CUDA-GDB sources are available as follows:

- For CUDA Toolkit 7.0 and newer, in the installation directory extras/. The directory is created by default during the toolkit installation unless the .rpm or .deb package installer is used. In this case, the cuda-gdb-src package must be manually installed.
- ► For CUDA Toolkit 6.5, 6.0, and 5.5, at https://github.com/NVIDIA/cuda-gdb.
- ► For CUDA Toolkit 5.0 and earlier, at ftp://download.nvidia.com/CUDAOpen64/.
- Upon request by sending an e-mail to mailto:oss-requests@nvidia.com.

# Chapter 2. CUDA 10.1 UPDATE 2 RELEASE NOTES

The release notes for the CUDA Toolkit can be found online at http://docs.nvidia.com/ cuda/cuda-toolkit-release-notes/index.html.

# 2.1. General CUDA

- CUDA 10.1 Update 2 is a minor update that is binary compatible with CUDA 10.1. This release will work with all versions of the R418 NVIDIA driver.
- All CUDA Libraries now adhere to independent semantic versioning. Because of this change, version numbers will differ between each of the libraries and the CUDA Toolkit.
- Support is added for RHEL 8 and EPEL.
- For RHEL 7, the <u>package manager installation</u> and the <u>package upgrade instructions</u> have been updated.

### 2.1.1. CUDA Tools

#### 2.1.1.1. CUDA Compilers

- The following compilers were added as host compilers in **nvcc**:
  - Xcode 10.2

# 2.2. CUDA Libraries

### 2.2.1. cuBLAS Library

 cuBLAS can now apply Tensor Core implementations for more number of problem sizes. See Tensor Core Restrictions.

### 2.2.2. cuSOLVER Library

 cuSOLVER has added support for multi-GPU execution, beginning in this release with a multi-GPU symmetric eigenvalue solver. Multi-GPU extensions are included in a separate binary/header with the Mg suffix, cuSOLVERMg.

### 2.2.3. cuSPARSE Library

Added the following features and enhancements for the cuSPARSE library:

- The cuSPARSE library now provides an API, cusparseConstrainedGEMM, to perform a dense matrix multiplication with sparse output constrained to a sparsity pattern set by the user.
- ► Generic cuSPARSE APIs **spvv**, **spMv**, and **spMM** are now supported on PowerPC Platforms.

### 2.2.4. NPP Library

- NPP now supports batched color conversion APIs (for packed-to-packed and planar-to-packed).
- NPP Geometry Transform API now supports FP16 format.

### 2.2.5. nvJPEG Library

- The nvJPEG library now supports transcoding.
- The nvJPEG encoder library now supports progressive encoding.
- Windows support is added for the nvJPEG library.

# 2.3. Deprecated Features

The following features are deprecated in this CUDA 10.1 Update 2 release. The features still work in the current release, but their documentation may have been removed, and they will become officially unsupported in a future release. We recommend that developers employ alternative solutions to these features in their software. **General CUDA** 

 CUDA 10.1 Update 2 and later versions are no longer supported on Ubuntu 18.10.

#### CUDA Tools

- Oracle JDK 8 JRE, required by Nsight Eclipse Edition and Visual Profiler, are no longer included in the CUDA Toolkit as of version 10.1 Update 2, due to Oracle upgrade licensing changes. For further details, see JDK references in:
  - ► Installing Nsight Eclipse Edition & Running Nsight Eclipse Edition
  - <u>Visual Profiler</u>

#### **CUDA** Libraries

Several cuSPARSE features have been deprecated and are planned for removal in CUDA 11. Users can see compile time warnings by defining the macro ENABLE\_CUSPARSE\_DEPRECATED to check compliance. The deprecated routines are noted in the documentation with their replacements where applicable.

### 2.4. Resolved Issues

#### 2.4.1. CUDA Tools

It is possible that in CUDA Toolkit 10.1 Update 1, the CUPTI tool loads libcuda.so instead of libcuda.so.1, leading to failures. This is fixed in CUDA 10.1 Update 2.

#### 2.4.2. CUDA Compilers

 It is possible that in CUDA Toolkit 10.1 Update 1, compiling the same PTX code with PTXAS produces incorrect results on sm\_60 architectures. This is fixed in CUDA 10.1 Update 2.

#### 2.4.3. CUDA Libraries

#### cuFFT Library

In <u>cuFFT library</u>, on sm\_60 and earlier architectures, for batched (batch > 1) 2D FFTs C2R with odd sizes, if the input is not Hermitian-conjugate, then the computation produced incorrect results. This is fixed in CUDA 10.1 Update 2.

#### cuSOLVER Library

Starting with CUDA 10.1 Update 2, the NVIDIA LAPACK library liblapack\_static.a is a subset of LAPACK, and only contains the GPUaccelerated stedc and bdsqr. The user must link libcusolver\_static.a with liblapack\_static.a in order to build the application successfully. See <u>"Link Thirdparty LAPACK Library."</u>

#### NPP Library

When nppiCFAToRGB\_8u\_C1C3R API is run in **NPPI\_BAYER\_GBRG** mode, it is possible that artifacts are generated in the images. This is fixed, so that now the artifacts are not generated.

#### nvJPEG Library

- Resolved the issue of ROI for non-split nvJPEG decoder API.
- Resolved the nvJPEG decoder issue of GPU\_HYBRID\_BACKEND requesting Y channel output for 3-channel input data.

#### cuBLAS Library

 Resolved the issue where Tensor Core implementations were not picked by heuristic even though they were applicable.

### 2.5. Known Issues

#### 2.5.1. General CUDA

For RHEL 7 the <u>package manager installation</u> and the <u>package upgrade instructions</u> have been updated. If you continue to use old instructions while installing CUDA 10.1 Update 2, then the installation will be broken. A symptom of this broken installation is that the entry "cuda-drivers" will not be listed in the list of packages to be installed.

## 2.5.2. CUDA Tools

#### **CUDA Compilers**

The maximum error bounds are revised in Tables 6 and 7 in <u>Section E: Mathematical</u> <u>Functions, Standard functions</u>, of the Programming Guide. The revised bounds are as follows:

#### Table 6:

normf: An error bound cannot be provided because a fast algorithm is used with accuracy loss due to round-off (revised from 4 ulp).

#### Table 7:

- sin, cos, sinpi, cospi, sincos, sincospi: 2 ulps (revised from 1 ulp)
- **acos**: 2 ulps (revised from 1 ulp)
- **erfc**: 5 (revised from 4)
- erfcx: 4 (revised from 3)
- normcdfinv: 8 (revised from 7)
- sinh: 2 (revised from 1)
- tanh: 2 (revised from 1)
- norm: An error bound cannot be provided because a fast algorithm is used with accuracy loss due to round-off (revised from 3 ulp)

Nsight Compute on Windows Subsystem for Linux is not supported.

Oracle JDK 8 JRE runtime, required by Nsight Eclipse Edition and Visual Profiler, is no longer being included in the CUDA Toolkit (as of version 10.1 Update 2), due to Oracle upgrade licensing changes. For further details, see JDK references in:

- Installing Nsight Eclipse Edition & Running Nsight Eclipse Edition
- <u>Visual Profiler</u>

For the Visual Profiler and Nsight Eclipse Edition, accessing the local help document leads to HTTP Error 500. The workaround is to access the NVIDIA online at these links below:

- Visual Profiler, and
- Nsight Eclipse Edition

#### **CUDA** Libraries

In <u>cuFFT library</u>, on sm\_60 and earlier architectures, for batched (batch > 1) 2D FFTs C2R with odd sizes, when the input is not Hermitian-conjugate, then the computation might be slower.

# Chapter 3. CUDA 10.1 UPDATE 1 RELEASE NOTES

The release notes for the CUDA Toolkit can be found online at http://docs.nvidia.com/ cuda/cuda-toolkit-release-notes/index.html.

# 3.1. General CUDA

- CUDA 10.1 Update 1 is a minor update that is binary compatible with CUDA 10.1. This release will work with all versions of the R418 NVIDIA driver.
- CUDA compatibility is now supported on R396 (>= 396.26). See documentation on using CUDA compatibility.
- Added support for Microsoft Visual Studio 2019 (RTM) including CUDA compiler and Nsight Visual Studio Edition integration.
- IBM POWER9 systems now support the NVIDIA Tesla T4 GPU.

# 3.2. CUDA Tools

The profiling APIs in the header cupti\_profiler\_target.h and the Perfworks metric APIs in the headers nvperf\_host.h and nvperf\_target.h are supported on the IBM POWER platform.

### 3.2.1. CUDA Compilers

- The following compilers are supported as host compilers in **nvcc**:
  - Clang 8.0
  - Microsoft Visual Studio 2019 (RTM and all updates)
- The tool nvcc now supports array capture in extended \_\_host\_\_\_\_device\_\_\_ and \_\_\_\_device\_\_\_ lambdas. See Programming Guide for Extended Lambda Restrictions.
- The representation of <u>host</u> <u>device</u> extended lambda in the code sent to the host compiler has been improved. As a result, the Clang host compiler is now able to inline the body of the <u>host</u> <u>device</u> extended lambda at the point of the call. This enhances the host side run-time performance for a <u>host</u>

\_\_device\_\_ extended lambda compared to previous CUDA versions when using the Clang host compiler.

# 3.3. CUDA Libraries

### 3.3.1. cuBLAS Library

This release features cuBLAS version 10.2.0 which adds the following features and enhancements to cuBLASLt API:

- cuBLASLt added the Tensor Core-accelerated IMMA kernels for Turing architecture GPUs with CUDA\_R\_8I outputs, float (and optionally per-row) scaling and optional ReLU and/or bias during epilogue phase of matrix multiplication.
- cuBLASLt extended the mixed precision Tensor Core-accelerated complex coverage to support both half (FP16) and single (FP32) precision outputs.
- cuBLASLt out-of-place reduction modes now work with arbitrary number of splits in the K dimension of the matrix multiply operation (i.e., split-K algorithm).
- Several performance improvements are made in this release. These enhancements are targeted for Tensor Core-accelerated mixed and single precision matrix multiplications on Volta and Turing GPU architectures.

### 3.3.2. cuSPARSE Library

Added the following features and enhancements for the cuSPARSE library:

- Added new generic routines for sparse vector-vector product (SpVV) and sparse matrix-vector multiplication (SpMV).
- Added new helper functions for sparse vector and matrix descriptors.
- Added support for CSR format in generic routines.
- Added new helper functions to handle return status and library version.

#### 3.3.3. cuFFT Library

- Added CUDA graphs support via stream capture. This support is limited to single GPU FFT plans.
- Improved the performance for most R2C/C2R FFTs.
- Multi-GPU dimension coverage extended for 2D/3D to not be limited by size decomposition by prime factors <= 127 alone for sizes up to 4096 in single precision (2048 in double precision).

#### 3.3.4. NPP Library

Added a new Label Marker API, nppiLabelMarkersUF\_xx\_C1R(), providing 40x-60x speedup for binary image label marker (4 or 8 connected component). This API is based on single pass GPU algorithm (Union Find).

### 3.3.5. nvJPEG Library

- The decoding of large images is sped up by adding GPU-based Huffman decode for single images.
- Introduced ROI decoding. This enables setting the coordinates for a user-defined region of interest (ROI) for decoding the image. This helps to decode any size and shape of a given image.
- Decoding now supports the color format CMYK.
- Introduced decoupled decoding. This enables the user to run the host phase of the decoding process independently of the device phase. Currently this is available only for single image multiphase decoding.

# 3.4. Deprecated Features

The following features are deprecated in the current release of the CUDA software. The features still work in the current release, but their documentation may have been removed, and they will become officially unsupported in a future release. We recommend that developers employ alternative solutions to these features in their software.

#### cuSPARSE Library

The routine cusparse<t>csrmm() is deprecated, and replaced by cusparseSpMM. See here.

# 3.5. Resolved Issues

### 3.5.1. General CUDA

- Fixed an issue with the CUDA .run installer where the .run installer would continue if the NVIDIA driver installation failed. Now the .run installer terminates if the driver installation fails for any reason.
- Improved the performance of **cudaMemGetInfo** by more than 40x (in some cases).
- Fixed a memory corruption issue in some cases when using Unified Memory prefetch hints on IBM POWER9 systems.

# 3.5.2. CUDA Compilers

- On PowerAI 1.5.4 with CUDA 10.0, attempting to use cuda-gdb might fail with: "internal-error: add\_cie: Assertion `n 1 || cie\_table-entries[n - 1]-cie\_pointer cie-cie\_pointer' failed." This is fixed in CUDA 10.1 Update 1.
- Fixed an issue with using the <u>\_\_\_attribute\_\_</u> on macOS with Xcode.
- Fixed possible compilation failures with using nvcc on macOS with Xcode 10. Nvcc now defaults to using C++03 as the default dialect on macOS.

- Fixed an error in nvcc where assignment operation in a template dependent context would fail to compile. For example, the below code will now compile successfully: reinterpret\_cast<float4\*>(&smem[0\*BYTES\_PER\_ROW\_WITH\_SKEW])[0] = sums\_[0];
- Fixed an issue with nvcc not correctly supporting type traits
   (<u>has\_trivial\_assign</u>) for integral types with Microsoft Visual Studio 2017.
- Fixed an issue with PTXAS when interpreting the debug information generated using LLVM (Clang 9 trunk).
- Fixed an issue in nvcc so that it will not generate warnings on unreachable code (to match other compilers such as gcc and Clang).
- Fixed an issue with nvcc when used with GCC 7.3 when building CUDA applications with member function calls.
- nvcc now supports the \_\_\_\_((deprecated)) variable attribute.
- Fixed an issue with parsing the \_\_integer\_pack intrinsic.
- Fixed an issue with variadic macro support in nvrtc. The below example code is now supported: #define \_fc\_vec\_elem(VEC, TYPE, ELEM, ...) (\*((TYPE \*)&(VEC) + (ELEM)))
- nvcc now handles the case where same value -std options can be specified on the command line.

### 3.5.3. CUDA Profiler

- Fixed a memory leak in the CUPTI tracing APIs.
- In CUDA 10.0 the command-line profiler nvprof overwrites LD\_PRELOAD with libaccinj64.so.10.1 when the option --openmp-profiling is enabled (which is the default), thereby preventing the use of software such as Spectrum MPI that relies on LD\_PRELOAD when using nvprof.
- Fixed an issue with nvprof to generate CSV for --analysis-metrics.

### 3.5.4. CUDA Libraries

#### NPP Library

Corrected floating point rounding error within the function **nppiGetResizeRect()** that for certain inputs resulted in slightly different destination image sizes (e.g., 999 pixels in a dimension where 1000 was expected).

#### cuSOLVER Library

- The new cusolverDn<t>getrf (LU factorization) delivers better performance and resolves potential race condition. Prior to CUDA 10.1 Update 1, getrf may hang on some GPUs.
- In the function cusolverSpXcsrqrAnalysis[host], if the zero fill-in exceeds the limits of 32-bit integer indexing, it leads to a segmentation fault. The function now returns CUSOLVER\_STATUS\_ALLOC\_FAILED if the zero fill-in exceeds the limits of 32-bit integer.
- In CUDA 10.1 if the user links their app with libcusolver\_static.a,
   liblapack\_static.a and 3rd party LAPACK library (for example, MKL), then
   [s]c]syevdx and [s]c]sygvdx may fail because some LAPACK symbols in

```
liblapack_static. a are replaced by MKL. This is fixed in CUDA 10.1 Update 1.
```

#### cuFFT Library

A memory corruption issue was fixed that affected FFTW API: **fftw\_execute\_dft** function, when using different pointers between FFTW planning and execution calls. **cuBLAS** 

The cuBLASLt documentation in CUDA 10.1 incorrectly stated that the mixed precision complex matrix multiply supported the single precision output, while it only supported the half precision output. This is corrected in CUDA 10.1 Update 1 by extending the support to both single and half precision outputs.

#### 3.5.5. CUDA Tools

 In CUDA 10.1 the tool nvprof overwrites LD\_PRELOAD with libaccinj64.so.10.1, thereby preventing the use of SMPI + nvprof with CUDA10.0/CUDA10.1. This is fixed in CUDA 10.1 Update 1.

### 3.6. Known Issues

#### 3.6.1. General CUDA

► The Tesla family GPUs on Windows systems will not be able to use profiling tools while using the driver provided by CUDA 10.1 Update 1 installers.

The Tesla Driver 425.25, provided with the CUDA Toolkit 10.1 Update 1, restricts access to the GPU performance counters used by the NVIDIA profiling tools (Nsight Compute, nvprof, Visual Profiler, CUPTI, and Nsight Visual Studio Edition).

Attempting to access the GPU performance counters will result in a reported error of **ERR\_NVGPUCTRPERM** and a link to https://developer.nvidia.com/ ERR\_NVGPUCTRPERM, for instructions on how system administrators can allow access to unprivileged users for profiling. On Windows systems using the Tesla family GPUs in TCC mode, the NVIDIA Control Panel should be used to grant access. However, the control panel in this specific configuration has a bug that prevents this process.

Follow any one of the three workarounds described below:

- 1. Run in WDDM mode. The control pane respects the settings in this mode.
- 2. Do not use the 425.25 Tesla driver shipped with the the CUDA Toolkit 10.1 Update 1.

Instead:

- a. Use the CUDA Toolkit 10.1 toolkit along with the 418.96 driver (released in Feb 2019).
- b. Use any of the following public drivers:
  - ► R418 (418.81)
  - ► R418 (418.91)

- ► R418 (419.17)
- ► R418 (419.35)
- Tesla driver 419.69
- Or any driver released after May 2019
- c. **DO NOT USE**: Any driver 419.67 or later unless it mentions a fix to this issue.
- 3. Launch the profiling tool with "Run as administrator."
- Power 9 + V100 performance is slower due to a hardware bug on some applications.

### 3.6.2. CUDA Tools

In Linux\_x64, Linux\_ppc64le, and Windows10 versions of CUDA 10.1, the CUPTI documentation might not have html or pdf copy. The CUPTI documentation is available online here.

#### 3.6.3. CUDA Libraries

 The cuSPARSE generic APIs are currently available only for Linux x86\_64 (AMD64) systems. Using these APIs on any other systems will result in compile-time or runtime failures.

# Chapter 4. CUDA 10.1 RELEASE NOTES

The release notes for the CUDA Toolkit can be found online at http://docs.nvidia.com/ cuda/cuda-toolkit-release-notes/index.html.

# 4.1. General CUDA

- ► Introducing NVIDIA<sup>®</sup> Nsight<sup>TM</sup> Systems, a system-wide performance analysis tool designed to visualize an application's algorithms. This tool will help you identify the largest opportunities to optimize, and efficiently tune to scale across any quantity or size of CPUs and GPUs—from a large server to the smallest SoC. See more here.
- Added 6.4 version of the Parallel Thread Execution instruction set architecture (ISA). For more details on new (noreturn, mma) and deprecated instructions (satfinite, non-sync versions of shfl and vote), see this section in the PTX documentation.
- The following new operating systems are supported by CUDA. See the System Requirements section in the NVIDIA CUDA Installation <u>Guide</u> for Linux for a full list of supported operating systems.
  - Ubuntu 18.10
  - ► RHEL 7.6
  - Fedora 29
  - ► SUSE SLES 12.4
  - Windows Server 2019
  - Windows 10 (October 2018 Update)
- Improved the scalability of cudaFree\* APIs on multi-GPU systems.
- Added support for cooperative group kernels (using the cudaLaunchCooperativeKernel API) with MPS.
- Relaxed IPC restrictions so that P2P can be enabled between devices that are not set by CUDA\_VISIBLE\_DEVICES.
- In CUDA 10.1 the <u>CUDA Runtime API error codes</u> are renumbered to match, wherever possible, their CUDA Driver API equivalents.
- Added GPU accounting, on Volta only, to keep track of open compute contexts and GPU utilization. This data is updated when the driver is loaded, but can be retrieved in driver-loaded and driver-unloaded modes via Out of Band (OOB).

- Added an out-of-band mechanism to fetch the instantaneous GPU and memory utilization See <u>SMBPBI spec</u> for the documentation.
- Added the ability to query GPU NVLink error rates / counts via out of band (OOB) with or without a driver present. See <u>SMBPBI spec</u> for the documentation.
- Added support for installing CUDA using runfiles on POWER (ppc64le) platforms.
- Added new CUDA samples for CUDA Graph APIs.

## 4.2. CUDA Tools

### 4.2.1. CUDA Compilers

- The following compilers are supported as host compilers in **nvcc**:
  - ► GCC 8.x
  - Clang 7.0
  - Microsoft Visual Studio 2017 (RTW, and all updates)
  - Microsoft Visual Studio 2019 (Preview releases)
  - ► PGI 19.x
  - ► ICC 19.0
  - Xcode 10.1 (10B61)
- New functions \_\_isShared(), \_\_isConstant() and \_\_isLocal() have been added, to check if a generic pointer points to an object in \_\_shared\_\_, \_\_constant\_\_ or local memory, respectively. These functions are documented in the CUDA C Programming Guide, along with the existing \_\_isGlobal() function.
- The existing API functions nvrtcGetLoweredName and nvrtcAddNameExpression have been enhanced to allow looking up the mangled (lowered) name of \_\_constant\_\_ and \_\_device\_\_ variables. Details and example here: <u>https://docs.nvidia.com/cuda/nvrtc/index.html#accessing-lowered-names</u>.
- nvcc now supports the "-MF" and "-MM" flags related to dependency generation. See below the description of the new flags from "nvcc --help":
  - --generate-nonsystem-dependencies (-MM) : Same as --generatedependencies but skip header files found in system directories (Linux only).
  - --dependency-output (-MF): Specify the output file for the dependency file generated with -M or -MM. If this option is not specified, the output is the same as if -E has been specified.

### 4.2.2. CUDA Profiler

- ► For new features in Visual Profiler and **nvprof**, see the <u>What's New</u> section in the Profiler User's Guide.
- For new features available in CUPTI, see the <u>What's New</u> section in the CUPTI documentation.
- For system wide profiling, use Nsight Systems. Refer to the Nsight Systems <u>Release</u> <u>Notes</u>.

 For profiling specific CUDA kernels, use Nsight Compute. Refer to the Nsight Compute <u>Release Notes</u>.

### 4.2.3. CUDA-MEMCHECK

• For new features in CUDA-MEMCHECK, see the Release Notes in the CUDA-MEMCHECK documentation.

# 4.3. CUDA Libraries

This release of the CUDA toolkit is packaged with libraries that deliver new and extended functionality, bug fixes, and performance improvements for single and multi-GPU environments.

Also in this release the **soname** of the libraries has been modified to not include the minor toolkit version number. For example, the cuFFT library **soname** has changed from **libcufft.so.10.1** to **libcufft.so.10**. This is done to facilitate any future library updates that do not include API breaking changes without the need to relink.

### 4.3.1. cuBLAS Library

- With this release, on Linux systems, the cuBLAS libraries listed below are now installed in the /usr/lib/<arch>-linux-gnu/ or /usr/lib64/ directories as shared and static libraries. Their interfaces are available in the /usr/include directory:
  - cublas (BLAS)
  - cublasLt (new Matrix Multiply library)
- Note that the new installation locations of cuBLAS libraries are different from the past versions. In the past versions the libraries were installed in directories under the main toolkit installation directory.
- Package managers on Linux OSs will remove the previous version of cuBLAS and update to the new libraries in the new location. For linking and execution make sure the new location is specified within your paths such as LD\_LIBRARY\_PATH.
- With this update, the versioning scheme of the cuBLAS library has changed to a 4-digit version. Because of this change, version numbers might differ between the CUDA toolkit and cuBLAS libraries in future releases. For the new 4-digit version:
  - The first three digits follow semantic versioning, and
  - The last digit is the build number.
- A new library, the cuBLASLt, is added. The cuBLASLt is a new lightweight library dedicated to GEneral Matrix-to-matrix Multiply (GEMM) operations with a new flexible API. This new library adds flexibility in matrix data layouts, input types, compute types, and also in choosing the algorithmic implementations and heuristics through parameter programmability. Read more at: <a href="http://docs.nvidia.com/cuda/cublas/index.html#using-the-cublasLt-api">http://docs.nvidia.com/cuda/cublas/index.html#using-the-cublasLt-api</a>.
- The new cuBLASLt library is packaged as a separate binary and a header file. Also, the cuBLASLt now adds support for:

- Utilization of IMMA tensor core operations on Turing GPUs for int8 input matrices.
- FP16 half-precision CGEMM split-complex matrix multiplies using tensor cores on Volta and Turing GPUs.

### 4.3.2. cuSOLVER Library

- For symmetric dense eigensolver:
  - Added a new selective eigensolver functionality for standard and generalized eigenvalue problems: SYEVDX and SYGVDX
  - Improved the performance for full eigenspectrum eigensolver.
- ► Added a new batched GESVDA API that computes the approximate singular value decomposition of a tall skinny **m×n** matrix A.
- Added a new POTRI API that computes the inverse of a symmetric positive definite matrix, using the Cholesky factorization computed by DPOTRF.

### 4.3.3. cuSPARSE Library

- Added a new generic Sparse x Dense Matrix Multiply (SpMM) APIs that encapsulates the functionality of many legacy APIs.
- Added a new COO matrix-matrix multiplication (cooMM) implementation with:
  - Deterministic and non-deterministic variants
  - Batched SpMM
  - Support for multiple data type combinations
  - Speed-ups w.r.t. csrMM for matrices with highly irregular nnzs/row
- Added two new algorithms for csr2csc format conversions with improved performance and reduced memory use.

### 4.3.4. cuFFT Library

- Improved the performance and scalability for the following use cases:
  - multi-GPU non-power of 2 transforms
  - R2C and Z2D odd sized transforms
  - 2D transforms with small sizes and large batch counts

### 4.3.5. cuRAND Library

- Improved the performance of the following random number generators:
  - MTGP32
  - MRG32k3a
  - Sobol32 and Scrambled Sobol32
  - Sobol64 and Scrambled Sobol64

### 4.3.6. NPP Library

- Some of the most commonly used image processing functions were extended to support the FP16 (<u>half</u>) data type on GPU architectures Volta and beyond.
- Added support for application-managed stream contexts. Application-managed stream contexts make NPP truely stateless internally, allowing for rapid, stream context switching with no overhead.
- While it is recommended that all new NPP application code use applicationmanaged stream contexts, existing application code can continue to use nppSetStream() and nppGetStream() to manage stream contexts (also with no overhead now). But over time NPP will likely deprecate the older non-applicationmanaged stream context API.

### 4.3.7. nvJPEG Library

- Added baseline encoding functionality to the library that will be extended in future releases.
- Added new batched decoding that uses GPU acceleration for all phases of computation. This delivers significant performance gains for large batches of images where most images are baseline encoded JPEG images.
- Added new APIs for pinned memory allocator and for memory overallocations.
- The nvJPEG library is now added to the Linux ppc64le CUDA Toolkit distributions.

# 4.4. Deprecated Features

The following features are deprecated in the current release of the CUDA software. The features still work in the current release, but their documentation may have been removed, and they will become officially unsupported in a future release. We recommend that developers employ alternative solutions to these features in their software.

#### **General CUDA**

- Nsight Eclipse Edition standalone is deprecated in CUDA 10.1, and will be dropped in the release that immediately follows CUDA 10.1.
- Support for RHEL 6.x is deprecated with CUDA 10.1. It may be dropped in a future release of CUDA. Customers are encouraged to adopt RHEL 7.x to use new versions of CUDA.
- The following compilers are no longer supported as host compilers for **nvcc** 
  - ► PGI 17.x
  - Microsoft Visual Studio 2010
  - Clang versions lower than 3.7
- Microsoft Visual Studio versions 2011, 2012 and 2013 are now deprecated as host compilers for nvcc. Support for these compilers may be removed in a future release of CUDA.
- 32-bit tools are no longer supported starting with CUDA 10.0.

- NVIDIA GPU Library Advisor (gpu-library-advisor) is now deprecated and will be removed in a future release of the toolkit.
- The non-sync definitions of warp shuffle functions (\_\_shfl, \_\_shfl\_up, \_\_shfl\_down, and \_\_shfl\_xor ) and warp vote functions (\_\_any, \_\_all, \_\_ballot) have been removed when compilation is targeting devices with compute capability 7.x and higher.
- For WMMA operations with floating point accumulators, the satf (saturateto-finite value) mode parameter is deprecated. Using it can lead to unexpected results. See http://docs.nvidia.com/cuda/cuda-c-programming-guide/ index.html#wmma-description for details.

#### **CUDA** Libraries

- The nvGRAPH library is deprecated. The library will no longer be shipped in future releases of the CUDA toolkit.
- The nppGetGpuComputeCapability function will be deprecated in the next NPP release. Users should instead call the cudaGetDevice() to get the GPU device ID, then call the function cudaDeviceGetAttribute() twice, once with the cudaDevAttrComputeCapabilityMajor parameter, and once with the cudaDevAttrComputeCapabilityMinor parameter.

# 4.5. Resolved Issues

### 4.5.1. CUDA Compilers

- In CUDA 9.2 nvprune crashes when running on a library with bss sections--for example, while pruning libcusparse\_static.a. This is fixed in CUDA 10.1.
- Within a template function, the CUDA compiler previously incorrectly allowed the use of an undeclared function. In CUDA 10.1, this compiler bug has been fixed and may cause diagnostics to be emitted for code that was previously incorrectly accepted. For example:

```
//-- template <typename T>
__global__ void foo(T in) { bar(in); }
int main() { foo<<<1,1>>>(1); }
__device__ void bar(int) { }
//--
```

This example was accepted by the previous CUDA compiler, but will generate a compile error with CUDA 10.1 compiler, because the function **bar()** is used in **foo()** before it has been declared.

### 4.5.2. CUDA Libraries

- In earlier releases, cuBLAS GEMM calls might randomly crash when running multiple host threads sharing one cuBLAS handle despite adhering to recommended usage in Thread Safety. This bug affects cuBLAS in earlier CUDA releases, and is fixed in CUDA 10.1.
- This bug affects cuSOLVER dense linear algebra functionality in CUDA 10.0, and is fixed in CUDA 10.1. An internal routine (LARFT) in cuSOLVER dense linear algebra

has a race condition if the user stream is not null. Although this routine is not in the cuSOLVER public API, it affects all routines based on householder reflection, including ORGBR, ORMTR, ORGTR, ORMQR, ORMQL, ORGQR AND SYEVD, SYGVD.

The following routines are **not** affected:

- ▶ GEQRF and GESVD.
- Also, GESVDJ and SYEVDJ are not affected because they are based-on the Jacobi method.
- The cuSOLVER routine SYEVD in CUDA 10.0 has a bug that may potentially impact single and single complex eigenvalue problems. In affected cases, the eigensolver may deliver inaccurate eigenvalues and eigenvectors. This bug only affects cuSOLVER dense linear algebra functionality in CUDA 10.0 and is fixed in CUDA 10.1.
- In CUDA 10.0 cuBLAS library, a bug in the batched LU factorization, cublas[S]
   D|C|Z]getrfBatched, may lead to wrong result or inconsistent result from run to run. CUDA 10.1 fixes the issue.

### 4.6. Known Issues

#### 4.6.1. General CUDA

On systems with a new install of Ubuntu 18.04.2, note that the installation of CUDA 10.1 and NVIDIA 418 drivers may result in the following error:

The following packages have unmet dependencies:

xserver-xorg-video-nvidia-418 : Depends:

```
xserver-xorg-core (>= 2:1.19.6-1ubuntu2~)
```

E: Unable to correct problems, you have held broken packages.

To recover from this error, install the *xserver-xorg-core* package and proceed with the installation of CUDA.

\$ sudo apt-get install xserver-xorg-core



This error is only observed on systems with a new install of Ubuntu 18.04.2. The error is not observed on systems that are upgraded from 18.04.1 to 18.04.2.

 Xwayland is not compatible with nvidia-settings and graphical CUDA samples. Recommend switching to Xorg session.

### 4.6.2. CUDA Tools

When using separate compilation and object linking for the device code (nvcc -relocatable-device-code=true or nvcc --device-c) the resulting binary may crash at runtime when all conditions below are true:

- 1. Some of the objects linked into the binary are generated by a previously released compiler (i.e., compiler from CUDA 10.0 Toolkit or earlier), and
- 2. Objects generated by the previously released compiler contain CUDA kernel function definition (i.e., "\_\_global\_\_" functions).

A possible workaround in such case is to generate all the objects with the CUDA 10.1 compiler.

- For known issues in cuda-memcheck, see the <u>Known Issues</u> section in the cudamemcheck documentation.
- For known issues in Nsight Compute, see the <u>Known Issues</u> section.
- When enabling the "Auto Profile" option in the Nsight Compute UI, profiling across different GPU architectures may fail.

To workaround this issue, profile the relevant kernels using the Nsight Compute CLI, or disable "Auto Profile" in the UI and manually profile these kernels.

- ► The tools **nv-nsight-cu** and **nv-nsight-cu-cli** will not work on any Linux platform with GLIBC version lower than 2.15. Hence these tools will not work on RHEL 6.10 and CentOS 6.10, which use GLIBC 2.12.
- For known issues in CUDA Profiling tools nvprof and Visual Profiler, see the <u>Profiler Known Issues</u> section in the Nsight Eclipse Edition Getting Started Guide.
- For known issues in the CUPTI library, see the <u>Limitations</u> section in the CUPTI document.
- For known issues in Nsight Eclipse, see the <u>Nsight Eclipse Known Issues</u> section in the Profiler User's Guide.
- On Windows CUPTI samples and other applications using the CUPTI APIs will result in the error "cupti.dll was not found". This is due to a mismatch in the CUPTI dynamic library name referenced in the import library "cupti.lib".

To workaround this issue rename the CUPTI dynamic library under the CUDA Toolkit directory (Default location is: "C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v10.1\extras\CUPTI\lib64") **from** "cupti64\_101.dll" **to** "cupti.dll".

► A call to cuptiFinalize()/cuptiProfilerDeInitialize() API can result in a hang with 410.x driver. Please use a 418 or later driver.

### 4.6.3. CUDA Libraries

 cuSOLVER dense linear algebra routine GETRF might exit with error code 702 on GPUs that have only 2 SMs Jetson TX1 GPUs.

# Chapter 5. THRUST V1.9.4 RELEASE NOTES

Thrust v1.9.4 adds asynchronous interfaces for parallel algorithms, a new allocator system including caching allocators and unified memory support, as well as a variety of other enhancements, mostly related to C++11/C++14/C++17/C++20 support.

The new asynchronous algorithms in the **thrust::async** namespace return **thrust::event** or **thrust::future** objects, which can be waited upon to synchronize with the completion of the parallel operation.

## 5.1. New Features

- thrust::event and thrust::future<T>, uniquely-owned asynchronous handles consisting of a state (ready or not ready), content (some value; for thrust::future only), and an optional set of objects that should be destroyed only when the future's value is ready and has been consumed.
  - The design is loosely based on C++11's std::future.
  - They can be .wait'd on, and the value of a future can be waited on and retrieved with .get or .extract.
  - Multiple thrust::events and thrust::futures can be combined with thrust::when\_all.
  - thrust::futures can be converted to thrust::events.
  - Currently, these primitives are only implemented for the CUDA backend and are C++11 only.
- New asynchronous algorithms that return thrust::event/thrust::futures, implemented as C++20 range style customization points:
  - thrust::async::reduce.
  - thrust::async::reduce\_into, which takes a target location to store the reduction result into.
  - thrust::async::copy, including a two-policy overload that allows explicit cross system copies which execution policy properties can be attached to.
  - thrust::async::transform.
  - thrust::async::for\_each.

- thrust::async::stable\_sort.
- thrust::async::sort.
- By default the asynchronous algorithms use the new caching allocators. Deallocation of temporary storage is deferred until the destruction of the returned thrust::future. The content of thrust::futures is stored in either device or universal memory and transferred to the host only upon request to prevent unnecessary data migration.
- Asynchronous algorithms are currently only implemented for the CUDA system and are C++11 only.
- exec.after(f, g, ...), a new execution policy method that takes a set of thrust::event/thrust::futures and returns an execution policy that operations on that execution policy should depend upon.
- New logic and mindset for the type requirements for cross-system sequence copies (currently only used by thrust::async::copy), based on:
  - thrust::is\_contiguous\_iterator and THRUST\_PROCLAIM\_CONTIGUOUS\_ITERATOR for detecting/indicating that an iterator points to contiguous storage.
  - thrust::is\_trivially\_relocatable and THRUST\_PROCLAIM\_TRIVIALLY\_RELOCATABLE for detecting/indicating that a type is memcpy-able (based on principles from <u>https://wg21.link/P1144</u>).
  - The new approach reduces buffering, increases performance, and increases correctness.
  - The fast path is now enabled when copying fp16 and CUDA vector types with thrust::async::copy.
- All Thrust synchronous algorithms for the CUDA backend now actually synchronize. Previously, any algorithm that did not allocate temporary storage (counterexample: thrust::sort) and did not have a computation-dependent result (counterexample: thrust::reduce) would actually be launched asynchronously.

Additionally, synchronous algorithms that allocated temporary storage would become asynchronous if a custom allocator was supplied that did not synchronize on allocation/deallocation, unlike cudaMalloc / cudaFree. So, now thrust::for\_each, thrust::transform, thrust::sort, etc are truly synchronous.

In some cases this may be a performance regression; if you need asynchrony, use the new asynchronous algorithms.

- Thrust's allocator framework has been rewritten. It now uses a memory resource system, similar to C++17's std::pmr but supporting static polymorphism. Memory resources are objects that allocate untyped storage and allocators are cheap handles to memory resources in this new model. The new facilities live in <thrust/mr/\*>.
  - thrust::mr::memory\_resource<Pointer>, the memory resource base class, which takes a (possibly tagged) pointer to void type as a parameter.
  - thrust::mr::allocator<T, MemoryResource>, an allocator backed by a memory resource object.
  - thrust::mr::polymorphic\_adaptor\_resource<Pointer>, a type-erased memory resource adaptor.

- thrust::mr::polymorphic\_allocator<T>, a C++17-style polymorphic allocator backed by a type-erased memory resource object.
- New tunable C++17-style caching memory resources, thrust::mr:: (disjoint\_)?(un)?synchronized\_pool\_resource, designed to cache both small object allocations and large repetitive temporary allocations. The disjoint variants use separate storage for management of the pool, which is necessary if the memory being allocated cannot be accessed on the host (e.g. device memory).
- System-specific allocators were rewritten to use the new memory resource framework.
- New thrust::device\_memory\_resource for allocating device memory.
- New thrust::universal\_memory\_resource for allocating memory that can be accessed from both the host and device (e.g. cudaMallocManaged).
- New thrust::universal\_host\_pinned\_memory\_resource for allocating memory that can be accessed from the host and the device but always resides in host memory (e.g. cudaMallocHost).
- thrust::get\_per\_device\_resource and thrust::per\_device\_allocator, which lazily create and retrieve a perdevice singleton memory resource.
- Rebinding mechanisms (rebind\_traits and rebind\_alloc) for thrust::allocator\_traits.
- thrust::device\_make\_unique, a factory function for creating a std::unique\_ptr to a newly allocated object in device memory.
- <thrust/detail/memory\_algorithms>, a C++11 implementation of the C+ +17 uninitialized memory algorithms.
- thrust::allocate\_unique and friends, based on the proposed C++23 std::allocate\_unique (https://wg21.link/P0211).
- New type traits and metaprogramming facilities. Type traits are slowly being migrated out of thrust::detail:: and <thrust/detail/\*>; their new home will be thrust:: and <thrust/type\_traits/\*>.
  - thrust::is\_execution\_policy.
  - thrust::is\_operator\_less\_or\_greater\_function\_object, which
    detects thrust::less, thrust::greater, std::less, and std::greater.
  - thrust::is\_operator\_plus\_function\_object, which detects thrust::plus and std::plus.
  - thrust::remove\_cvref(\_t)?, a C++11 implementation of C++20's
    thrust::remove\_cvref(\_t)?.
  - thrust::void\_t, and various other new type traits.
  - thrust::integer\_sequence and friends, a C++11 implementation of C++20's std::integer\_sequence.
  - thrust::conjunction, thrust::disjunction, and thrust::disjunction, a C++11 implementation of C++17's logical metafunctions.
  - Some Thrust type traits (such as thrust::is\_constructible) have been redefined in terms of C++11's type traits when they are available.
- <thrust/detail/tuple\_algorithms.h>, new std::tuple algorithms:

- thrust::tuple\_transform.
- thrust::tuple\_for\_each.
- thrust::tuple\_subset.
- Miscellaneous new std::-like facilities:
  - thrust::optional, a C++11 implementation of C++17's std::optional.
  - **thrust:** : **addressof**, an implementation of C++11's **std:** : **addressof**.
  - thrust::next and thrust::prev, an implementation of C++11's std::next and std::prev.
  - thrust::square, a <functional> style unary function object that multiplies its argument by itself.
  - <thrust/limits.h> and thrust::numeric\_limits, a customized version
     of <limits> and std::numeric\_limits.
- <thrust/detail/preprocessor.h>, new general purpose preprocessor facilities:
  - **THRUST PP CAT[2-5]**, concatenates two to five tokens.
  - THRUST\_PP\_EXPAND (\_ARGS) ?, performs double expansion.
  - **THRUST\_PP\_ARITY** and **THRUST\_PP\_DISPATCH**, tools for macro overloading.
  - **THRUST\_PP\_BOOL**, boolean conversion.
  - **THRUST\_PP\_INC** and **THRUST\_PP\_DEC**, increment/decrement.
  - **THRUST\_PP\_HEAD**, a variadic macro that expands to the first argument.
  - THRUST\_PP\_TAIL, a variadic macro that expands to all its arguments after the first.
  - **THRUST\_PP\_IIF**, bitwise conditional.
  - THRUST\_PP\_COMMA\_IF, and THRUST\_PP\_HAS\_COMMA, facilities for adding and detecting comma tokens.
  - THRUST\_PP\_IS\_VARIADIC\_NULLARY, returns true if called with a nullary \_\_VA\_ARGS\_\_.
  - **THRUST\_CURRENT\_FUNCTION**, expands to the name of the current function.
- New C++11 compatibility macros:
  - THRUST\_NODISCARD, expands to [[nodiscard]] when available and the best equivalent otherwise.
  - THRUST\_CONSTEXPR, expands to constexpr when available and the best equivalent otherwise.
  - THRUST\_OVERRIDE, expands to override when available and the best equivalent otherwise.
  - THRUST\_DEFAULT, expands to = default; when available and the best equivalent otherwise.
  - THRUST\_NOEXCEPT, expands to noexcept when available and the best equivalent otherwise.
  - THRUST\_FINAL, expands to final when available and the best equivalent otherwise.
  - THRUST\_INLINE\_CONSTANT, expands to inline constexpr when available and the best equivalent otherwise.
- <thrust/detail/type\_deduction.h>, new C++11-only type deduction helpers:

- **THRUST\_DECLTYPE\_RETURNS**\*, expand to function definitions with suitable conditional noexcept qualifiers and trailing return types.
- THRUST\_FWD(x), expands to ::std::forward<decltype(x)>(x).
- **THRUST\_MVCAP**, expands to a lambda move capture.
- THRUST\_RETOF, expands to a decltype computing the return type of an invocable.

## 5.2. New Examples

mr\_basic demonstrates how to use the new memory resource allocator system.

# 5.3. Other Enhancements

### 5.3.1. Tagged Pointer Enhancements

- New thrust::pointer\_traits specialization for void const\*.
  - nullptr support to Thrust tagged pointers.
  - New explicit operator bool for Thrust tagged pointers when using C++11 for std::unique\_ptr interoperability.
  - Added thrust::reinterpret\_pointer\_cast and thrust::static\_pointer\_cast for casting Thrust tagged pointers.

#### 5.3.2. Iterator Enhancements

- thrust::iterator\_system is now SFINAE friendly.
  - Removed cv qualifiers from iterator types when using thrust::iterator\_system.
  - Static assert enhancements:
  - New THRUST\_STATIC\_ASSERT\_MSG, takes an optional string constant to be used as the error message when possible.
  - Update THRUST\_STATIC\_ASSERT (\_MSG) to use C++11's static\_assert when it's available.
  - Introduce a way to test for static assertions.

### 5.3.3. Testing Enhancements

- Additional scalar and sequence types, including non-builtin types and vectors with unified memory allocators, have been added to the list of types used by generic unit tests.
- The generation of random input data has been improved to increase the range of values used and catch more corner cases.
- New truncate\_to\_max\_representable utility for avoiding the generation of ranges that cannot be represented by the underlying element type in generic unit test code.

- The test driver now synchronizes with CUDA devices and check for errors after each test, when switching devices, and after each raw kernel launch.
- The **warningtester** uber header is now compiled with NVCC to avoid needing to disable CUDA-specific code with the preprocessor.
- Fixed the unit test framework's ASSERT\_\* to print chars as ints.
- New **DECLARE\_INTEGRAL\_VARIABLE\_UNITTEST** test declaration macro.
- New DECLARE\_VARIABLE\_UNITTEST\_WITH\_TYPES\_AND\_NAME test declaration macro.
- thrust::system\_error in the CUDA backend now print out its cudaError\_t enumerator in addition to the diagnostic message.
- Stopped using conditionally signed types like **char**.

## 5.4. Resolved Issues

- Fixed compilation error when using <u>\_\_\_\_\_\_\_lambdas</u> with reduce on MSVC.
- Static asserted that thrust::generate / thrust::fill doesn't operate on const iterators.
- Fixed compilation failure with thrust::zip\_iterator and thrust::complex<float>.
- Fixed dispatch for the CUDA backend's thrust::reduce to use two functions (one with the pragma for disabling exec checks, one with THRUST\_RUNTIME\_FUNCTION) instead of one. This fixes a regression with device compilation that started in CUDA 9.2.
- Added missing <u>host</u> <u>device</u> annotations to a thrust::complex::operator= to satisfy GoUDA.
- Made thrust::vector\_base::clear not depend on the element type being default constructible.
- Removed flaky simple\_cuda\_streams example.
- Added missing thrust::device\_vector constructor that takes an allocator parameter.
- Updated the **range\_view** example to not use device-side launch.
- Ensured that sized unit tests that use **counting\_iterator** perform proper truncation.
- Refactored questionable copy\_if unit tests.

# Chapter 6. CUDA TEGRA RELEASE NOTES

The release notes for CUDA Tegra contain only information that is specific to the following:

- CUDA Tegra Driver, and
- Mobile version of other CUDA components such as: compilers, tools, libraries, and samples.

The release notes for the desktop version of CUDA also apply to CUDA Tegra. On Tegra, the CUDA Toolkit version is 10.1.

### 6.1. New Features

#### **CUDA Tegra Driver**

- Support is added for GPUDirect RDMA on AGX Jetson platform. This now enables a direct path for data exchange between the GPU and third-party peer devices using standard features of PCI Express.
- Support for Android P added.
- Error Reporting enriched in CUDA on mobile RM.
- Support added for Ubuntu 18.04 for the AGX Drive platform.
- Resumed the support for Ubuntu 16.04 host on the AGX Jetson platform.
- Performance optimizations that were previously enabled for QNX are now also available on Linux through user-mode submits.

#### **CUDA Compiler**

- Support added for GCC 7.3 on AGX Jetson platforms.
- Support added for CLANG 6.0.2 for NVCC on Android platforms.

# 6.2. Known Issues and Limitations

#### CUDA Tegra Driver

- Only the below color formats are supported for Vulkan-CUDA interoperability on Jetson and Android:
  - ► VK\_FORMAT\_R8\_UNORM
  - ► VK\_FORMAT\_R8\_SNORM
  - ► VK\_FORMAT\_R8\_UINT
  - ► VK\_FORMAT\_R8\_SINT
  - ► VK\_FORMAT\_R8\_SRGB
  - ► VK\_FORMAT\_R8G8\_UNORM
  - ► VK\_FORMAT\_R8G8\_SNORM
  - ► VK\_FORMAT\_R8G8\_UINT
  - ► VK FORMAT R8G8 SINT
  - ► VK\_FORMAT\_R8G8\_SRGB
  - ► VK\_FORMAT\_R16\_UNORM
  - ► VK\_FORMAT\_R16\_SNORM
  - ► VK FORMAT R16 UINT
  - ► VK\_FORMAT\_R16\_SINT
  - ► VK\_FORMAT\_R16\_SFLOAT

Other formats are currently not supported.

#### **CUDA Tools**

On NVIDIA DRIVE OS Linux systems, when using Nsight Compute CLI with "--mode attach" to attach to another process on the same machine, "--hostname 127.0.0.1" must be passed. This is because the default value of "localhost" for the "--hostname" parameter does not work.

# 6.3. Resolved Issues

#### **General CUDA**

- **CUDA-GDB on Linux:** The **set cuda memcheck on** command in cuda-gdb does not have any effect. This is fixed in CUDA 10.1.
- CUDA-GDB on QNX: ntoaarch64-gdb and cuda-qnx-gdb may hang when executing the run command. This is fixed in CUDA 10.1.

# 6.4. Deprecated Issues

#### **General CUDA**

• Deprecating support for Pascal product on Android.

#### Acknowledgments

NVIDIA extends thanks to Professor Mike Giles of Oxford University for providing the initial code for the optimized version of the device implementation of the double-precision exp() function found in this release of the CUDA toolkit.

NVIDIA acknowledges Scott Gray for his work on small-tile GEMM kernels for Pascal. These kernels were originally developed for OpenAI and included since cuBLAS 8.0.61.2.

#### 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.

#### 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



 $^{\odot}$  2007-2019 NVIDIA Corporation. All rights reserved. www.nvidia.com