



DATACENTER GPU MANAGER 2.0

v2.0 | July 2020

Release Notes



TABLE OF CONTENTS

- Changelog..... iii
- Patch Releases..... iii
 - DCGM v2.0.15..... iii
 - DCGM v2.0.13..... iii
- DCGM v2.0 GA..... iv
 - New Features..... iv
 - Improvements..... iv
 - Bug Fixes..... v
 - Known Issues..... v

CHANGELOG

This version of **DCGM** (v2.0) requires a minimum R418 driver that can be downloaded from **NVIDIA Drivers**. On NVSwitch based systems such as DGX A100 or HGX A100, a minimum of Linux R450 ($\geq 450.80.02$) driver is required. If using the new profiling metrics capabilities in DCGM, then a minimum of Linux R418 ($\geq 418.87.01$) driver is required. It is recommended to install the latest datacenter driver from NVIDIA drivers downloads site for use with DCGM.

Patch Releases

DCGM v2.0.15

DCGM v2.0.15 released in January 2021.

Improvements

- ▶ Added support for the NVIDIA A100 80GB product.
- ▶ Added new return codes to distinguish between warnings and errors when running DCGM Diagnostics.
- ▶ Changed the output for `dcgmi discovery -c` to list the available GPU Instances, Compute Instances and GPU UUIDs for A100 MIG.
- ▶ Added license files into the DCGM installer packages.

Bug Fixes

- ▶ Fixed an issue where DCGM Diagnostics would report a failed NVLink health status when running in MIG mode for A100 (since NVLink is not supported when in MIG mode).
- ▶ Fixed an issue where some profiling metrics are reported as **N/A** rather than **0s** when metrics monitoring is started.
- ▶ Fixed a bug in the PCIe checks in DCGM Diagnostics that would result in a crash in some cases.

DCGM v2.0.13

DCGM v2.0.13 released in October 2020.

Improvements

- ▶ Added support for attributing GPU telemetry to MIG devices on A100.
- ▶ Added support for Arm64 server platforms.

Bug Fixes

- ▶ Added an error message to indicate failure of DCGM Diagnostics on NVSwitch systems when Fabric Manager was not running. The test would previously fail with a `cuInit` error.

DCGM v2.0 GA

DCGM v2.0.10 released in July 2020.

New Features

General

- ▶ Added support for NVIDIA A100 (GPUs and NVSwitch based systems such as DGX A100 and HGX A100)
- ▶ Added support for NVIDIA A100 Multi-Instance GPU (MIG):
 - ▶ DCGM can enumerate GPU Instances (I) and GPU Compute Instances (CI)
 - ▶ Added the ability to monitor GPU-Is and GPU-CIs
- ▶ Added support for new A100 SKUs to the DCGM GPU Diagnostics
- ▶ DCGM 2.0 no longer includes the Fabric Manager (FM) for NVSwitch systems. FM is a separate package that needs to be installed with the R450 driver. DCGM 2.0 cannot be used on NVSwitch systems (e.g. DGX or HGX) that are running driver versions < R450.
- ▶ Added the ability (`dcgmHealthSet_v2` API) to set update interval and quota policy for health checks.
- ▶ Added support for CUDA 11 to DCGM GPU Diagnostics

Improvements

General

- ▶ DCGM 2.0.10 has lowered the minimum `glibc` requirement to 2.12 instead of 2.17.
- ▶ DCGM logs are no longer encrypted.
- ▶ The DCGM network protocol has been updated for performance and security. You cannot connect a 1.7.x DCGM library (`libdcgm.so`) to a 2.0.x `nv-hostengine` or vice versa. This includes `dcgmi` and using APIs like `dcgmConnect`.
- ▶ DCGM now supports 32 GPUs in a system (up from 16) (see `DCGM_MAX_NUM_DEVICES`).
- ▶ Updated APIs to support 3rd generation NVLink (`DCGM_NVLINK_MAX_LINKS_PER_GPU`) to 12 links per GPU.

- ▶ DCGM documentation can now be found online at <http://docs.nvidia.com/datacenter/dcgm> and packages no longer include documentation.

Bug Fixes

- ▶ Fixed an issue with `dcgmi` which could result in a crash when an invalid GPU list is provided via the `-i` option
- ▶ Fixed an issue with excessive CPU overhead when using the `dcgmHealthCheck` with `DCGM_HEALTH_WATCH_MEM`
- ▶ Fixed an issue where using profiling metrics with T4 in GPU VM passthrough, DCGM may report memory bandwidth utilization to be 12% higher.
- ▶ Fixed an issue where using multiplexing of profiling metrics, the PCIe bandwidth numbers returned by DCGM may be incorrect.

Known Issues

- ▶ On DGX-2/HGX-2 systems, ensure that `nv-hostengine` and the Fabric Manager service are started before using `dcgmproftester` for testing the new profiling metrics. See the Getting Started section in the DCGM User Guide for details on installation.
- ▶ On K80s, `nvidia-smi` may report hardware throttling (`clocks_throttle_reasons.hw_slowdown = ACTIVE`) during DCGM Diagnostics (Level 3). The stressful workload results in power transients that engage the HW slowdown mechanism to ensure that the Tesla K80 product operates within the power capping limit for both long term and short term timescales. For Volta or later Tesla products, this reporting issue has been fixed and the workload transients are no longer flagged as "HW Slowdown". The NVIDIA driver will accurately detect if the slowdown event is due to thermal thresholds being exceeded or external power brake event. It is recommended that customers ignore this failure mode on Tesla K80 if the GPU temperature is within specification.
- ▶ To report NVLINK bandwidth utilization DCGM programs counters in the HW to extract the desired information. It is currently possible for certain other tools a user might run, including `nvprof`, to change these settings after DCGM monitoring begins. In such a situation DCGM may subsequently return errors or invalid values for the NVLINK metrics. There is currently no way within DCGM to prevent other tools from modifying this shared configuration. Once the interfering tool is done a user of DCGM can repair the reporting by running `nvidia-smi nvlink -sc 0bz`; `nvidia-smi nvlink -sc 1bz`.

Notice

THE INFORMATION IN THIS GUIDE AND ALL OTHER INFORMATION CONTAINED IN NVIDIA DOCUMENTATION REFERENCED IN THIS GUIDE IS PROVIDED “AS IS.” NVIDIA MAKES NO WARRANTIES, EXPRESSED, IMPLIED, STATUTORY, OR OTHERWISE WITH RESPECT TO THE INFORMATION FOR THE PRODUCT, AND EXPRESSLY DISCLAIMS ALL IMPLIED WARRANTIES OF NONINFRINGEMENT, MERCHANTABILITY, AND FITNESS FOR A PARTICULAR PURPOSE. Notwithstanding any damages that customer might incur for any reason whatsoever, NVIDIA’s aggregate and cumulative liability towards customer for the product described in this guide shall be limited in accordance with the NVIDIA terms and conditions of sale for the product.

THE NVIDIA PRODUCT DESCRIBED IN THIS GUIDE IS NOT FAULT TOLERANT AND IS NOT DESIGNED, MANUFACTURED OR INTENDED FOR USE IN CONNECTION WITH THE DESIGN, CONSTRUCTION, MAINTENANCE, AND/OR OPERATION OF ANY SYSTEM WHERE THE USE OR A FAILURE OF SUCH SYSTEM COULD RESULT IN A SITUATION THAT THREATENS THE SAFETY OF HUMAN LIFE OR SEVERE PHYSICAL HARM OR PROPERTY DAMAGE (INCLUDING, FOR EXAMPLE, USE IN CONNECTION WITH ANY NUCLEAR, AVIONICS, LIFE SUPPORT OR OTHER LIFE CRITICAL APPLICATION). NVIDIA EXPRESSLY DISCLAIMS ANY EXPRESS OR IMPLIED WARRANTY OF FITNESS FOR SUCH HIGH RISK USES. NVIDIA SHALL NOT BE LIABLE TO CUSTOMER OR ANY THIRD PARTY, IN WHOLE OR IN PART, FOR ANY CLAIMS OR DAMAGES ARISING FROM SUCH HIGH RISK USES.

NVIDIA makes no representation or warranty that the product described in this guide will be suitable for any specified use without further testing or modification. Testing of all parameters of each product is not necessarily performed by NVIDIA. It is customer’s sole responsibility to ensure the product is suitable and fit for the application planned by customer and to do the necessary testing for the application in order to avoid a default of the application or the product. Weaknesses in customer’s product designs may affect the quality and reliability of the NVIDIA product and may result in additional or different conditions and/or requirements beyond those contained in this guide. NVIDIA does not accept any liability related to any default, damage, costs or problem which may be based on or attributable to: (i) the use of the NVIDIA product in any manner that is contrary to this guide, or (ii) customer product designs.

Other than the right for customer to use the information in this guide with the product, no other license, either expressed or implied, is hereby granted by NVIDIA under this guide. Reproduction of information in this guide is permissible only if reproduction is approved by NVIDIA in writing, is reproduced without alteration, and is accompanied by all associated conditions, limitations, and notices.

Trademarks

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

Copyright

© 2013-2021 NVIDIA Corporation. All rights reserved.