



NVIDIA L4 GPU Accelerator

Product Brief

Document History

PB-11316-001_v01

Version	Date	Authors	Description of Change
01	March 9, 2023	AV, SM	Initial release

Table of Contents

- Overview..... 1
- Specifications 2
 - Product Specifications 2
 - Environmental and Reliability Specifications..... 4
- Airflow Direction Support..... 5
- Product Features 6
 - PCI Express Interface Specifications 6
 - PCIe Support..... 6
 - Single Root I/O Virtualization Support..... 6
 - Interrupt Messaging..... 6
 - Polarity Inversion and Lane Reversal Support 7
 - Root of Trust..... 7
 - Form Factor 7
 - Hockey Stick Board Retention 9
- Support Information..... 10
 - Certifications..... 10
 - Agencies 10
 - Languages 11

List of Figures

Figure 1. NVIDIA L4 NVFF 5.5 HHL with Full Height Bracket.....	1
Figure 2. NVIDIA L4 Airflow Direction	5
Figure 3. NVIDIA L4 PCIe Card Dimensions with Full Height Bracket.....	8
Figure 4. NVIDIA L4 PCIe Card Dimensions with Low Profile Bracket.....	8
Figure 5. NVIDIA L4 Hockey Stick Tab.....	9

List of Tables

Table 1. Product Specifications	2
Table 2. Memory Specifications.....	3
Table 3. Software Specifications.....	3
Table 4. Board Environmental and Reliability Specifications	4
Table 5. Languages Supported.....	11

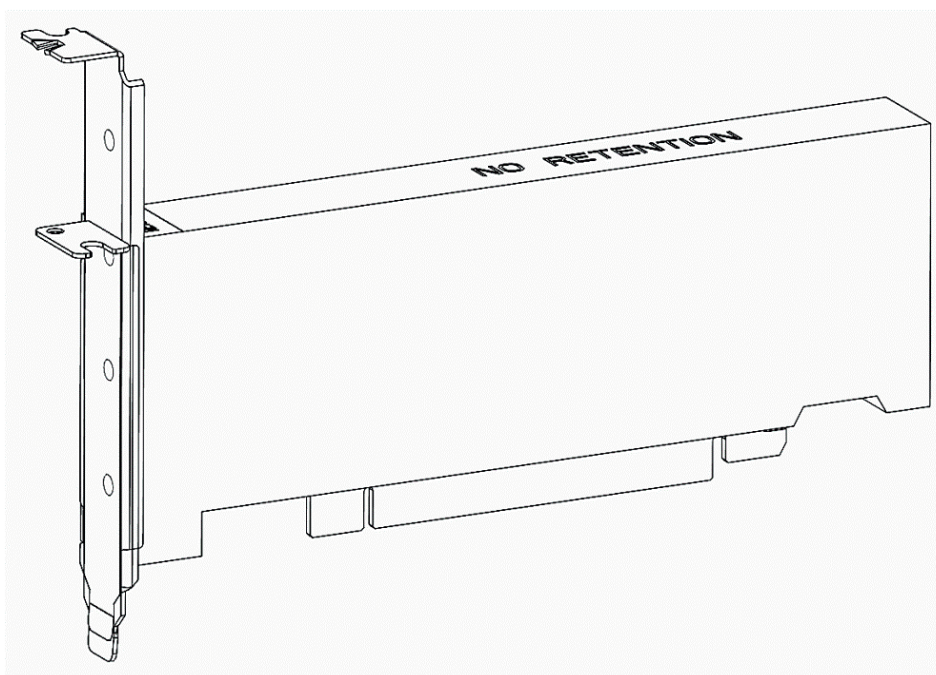
Overview

The NVIDIA L4 Tensor Core GPU delivers a versatile platform to accelerate Deep Learning, Graphics and Video processing applications in the Cloud and at the Edge. It is a half-height (low profile), half-length, single slot card featuring 24 GB of GDDR6 memory, x16 PCIe Gen4 connectivity at a 72 W maximum power envelope. It is a passively cooled card with a superior thermal design-requiring system airflow to operate and handles challenging ambient environments with ease (NEBS-3 capable).

Powered by the NVIDIA Ada Lovelace architecture, L4 provides revolutionary multi-precision performance to accelerate deep learning and machine learning training and inference, video transcoding, AI audio (AU) and video effects, rendering, data analytics, virtual workstations, virtual desktop, and many other workloads.

As part of NVIDIA AI, the L4 supports all AI frameworks and neural network models, delivering dramatic performance and efficiency that maximizes the utility of at-scale deployments.

Figure 1. NVIDIA L4 NVFF 5.5 HHHL with Full Height Bracket



Specifications

Product Specifications

Table 1 through Table 3 the product, memory, and software specifications for the NVIDIA L4 PCIe card.

Table 1. Product Specifications

Specification	NVIDIA L4
Product SKU	PG 193 SKU 200 NVPN: 699-2G193-0200-xxx
Total board power	72 W default 72 W maximum 40 W minimum
Thermal solution	Passive
Mechanical form factor	HHHL-SS (half-height, half-length, single-slot)
PCI Device IDs	Device ID: 0x27B8 Vendor ID: 0x10DE Sub-Vendor ID: 0x10DE Sub-System ID: 0x16CA
Four-part ID (VID:DEVID:SVID:SSID) ¹	10DE:27B8:10DE:16CA
GPU clocks	Base: 795 MHz Boost: 2,040 MHz
VBIOS	EEPROM size: 16 Mbit UEFI: Supported
Drivers	Linux: R525 or later Windows: R525 or later
PCI Express interface	Physical x16 PCIe lanes PCIe Gen4 x16, x8; Gen3 x16 Lane and polarity reversal supported
Performance states	P0, P8
Zero Power	Not supported

Specification	NVIDIA L4
Weight	Board: 270 Grams (excluding bracket) Bracket (Full height) with screws: 14 Grams Bracket (Half height) with screws: 9 Grams
Note: ¹ The NVIDIA L4 is uniquely identified by its complete four-part ID.	

Table 2. Memory Specifications

Specification	Description
Memory clock	6,251 MHz
Memory type	GDDR6
Memory size	24 GB
Memory bus width	192 bits
Peak memory bandwidth	300 GB/sec

Table 3. Software Specifications

Specification	NVIDIA L4
SR-IOV support	Supported: 32 VF (virtual functions)
BAR address (physical function)	BAR0: 16 MiB ¹ BAR1: 32 GiB ¹ BAR3: 32 MiB ¹
BAR address (virtual function)	BAR0: 8 MiB (256 KiB per VF) ¹ BAR1: 64 GiB, 64-bit (2 GiB per VF) ¹ BAR3: 1 GiB, 64-bit (32 MiB per VF) ¹
Message signaled interrupts	MSI-X: Supported MSI: Not supported
ARI Forwarding	Supported
Secure Boot	Supported (See “Root of Trust” section)
NVIDIA® CUDA® support	CUDA 12.0 or later
Virtual GPU software support	Supports vGPU 15.2 or later
PCI class code	0x03 – Display controller
PCI sub-class code	0x02 – 3D controller
ECC support	Enabled (by default); can be disabled using software
SMBus (8-bit address)	0x9E (write), 0x9F (read)
IPMI FRU EEPROM I2C address	0x50 (7-bit), 0xA0 (8-bit)
Reserved I2C addresses	0xA0, 0xAA, 0xAC
SMBus direct access	Supported

Specification	NVIDIA L4
SMBPBI (SMBus Post-Box Interface)	Supported
Note: ¹ The KiB, MiB, and GiB notations emphasize the “power of two” nature of the values. Thus, > 256 KiB = 256 × 1024 > 16 MiB = 16 × 1024 ² > 64 GiB = 64 × 1024 ³	

Environmental and Reliability Specifications

Table 4 provides the environment conditions specifications for the L4 PCIe card.

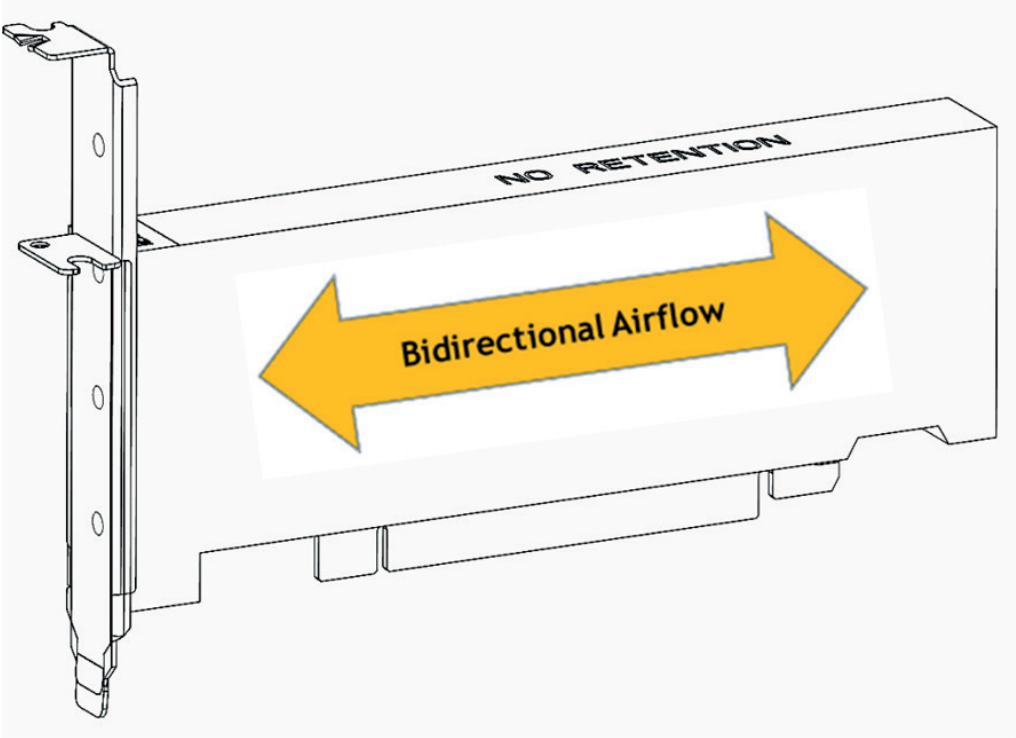
Table 4. Board Environmental and Reliability Specifications

Specification	Description
Ambient operating temperature	0°C to 50°C
Ambient operating temperature (short term) ¹	-5°C to 55°C
Storage temperature	-40°C to 75°C
Operating humidity (short term) ¹	5% to 93% relative humidity
Operating humidity	5% to 85% relative humidity
Storage humidity	5% to 95% relative humidity
Mean time between failures (MTBF)	Uncontrolled environment: ² 2,147,604 hours at 35°C Controlled environment: ³ 2,785,669 hours at 35°C
Notes: Specifications in this table are applicable up to 6,000 feet. ¹ A period not more than 96 hours consecutive, not to exceed 15 days per year. ² Some environmental stress with limited maintenance (GF35). ³ No environmental stress with optimum operation and maintenance (GB35).	

Airflow Direction Support

The NVIDIA L4 PCIe card employs a bidirectional heat sink, which accepts airflow either left-to-right or right-to-left directions.

Figure 2. NVIDIA L4 Airflow Direction



Product Features

PCI Express Interface Specifications

The following subsections describe the PCIe interface specifications for the L4 PCIe card.

PCIe Support

The L4 card supports PCIe Gen4. Gen4 x16 interface should be used when connecting to the L4 PCIe card.

Single Root I/O Virtualization Support

Single Root I/O (SR-IOV) Virtualization is a PCIe specification that allows a physical PCIe device to appear as multiple physical PCIe devices. Per PCIe specification, each device can have up to a maximum of 256 virtual functions (VFs). The actual number can depend on the device. SR-IOV is enabled in L4 PCIe card. The number of VFs supported is given in Table 3.

For each device, SR-IOV identifies two function classes:

- > Physical functions (PFs) constitute full-featured functionality. They are fully configurable, and their configuration can control the entire device. Naturally, a PF also has full ability to move data in and out of the device.
- > Virtual functions (VFs), which lack configuration resources. VFs exist on an underlying PF, which may support many such VFs. VFs can only move data in and out of the device; they cannot be configured and cannot be treated like a full PCIe device. The OS or hypervisor instance must be aware that they are not full PCIe devices.

The L4 requires that SBIOS and software support in the OS instance or hypervisor is configured to enable support SR-IOV. The OS instance or hypervisor must be able to detect and initialize PFs and VFs.

Interrupt Messaging

The L4 PCIe card only supports the MSI-X interrupt messaging protocol. The MSI interrupt protocol is not supported.

Polarity Inversion and Lane Reversal Support

Lane Polarity Inversion, as defined in the PCIe specification, is supported on the L4 PCIe card.

Lane Reversal, as defined in the PCIe specification, is supported on the L4 PCIe card. When reversing the order of the PCIe lanes, the order of both the Rx lanes and the Tx lanes must be reversed.

Root of Trust

The NVIDIA L4 provides a primary root of trust within the GPU that provides the following:

- > Secure boot
- > Secure firmware upgrade
- > Firmware rollback protection
- > Support for in-band firmware update disable (established after each GPU reset)
- > Secure application processor recovery

Form Factor

The NVIDIA L4 PCIe card conforms to NVIDIA Form Factor 5.5 specification for a half-height (low profile) half-length (HHHL) single slot PCIe card. For details refer to the *NVIDIA Form Factor 5.5 for Enterprise PCIe Products Specification* (NVOnline reference number 106337).

In this product brief, nominal dimensions are shown.

Figure 3. NVIDIA L4 PCIe Card Dimensions with Full Height Bracket

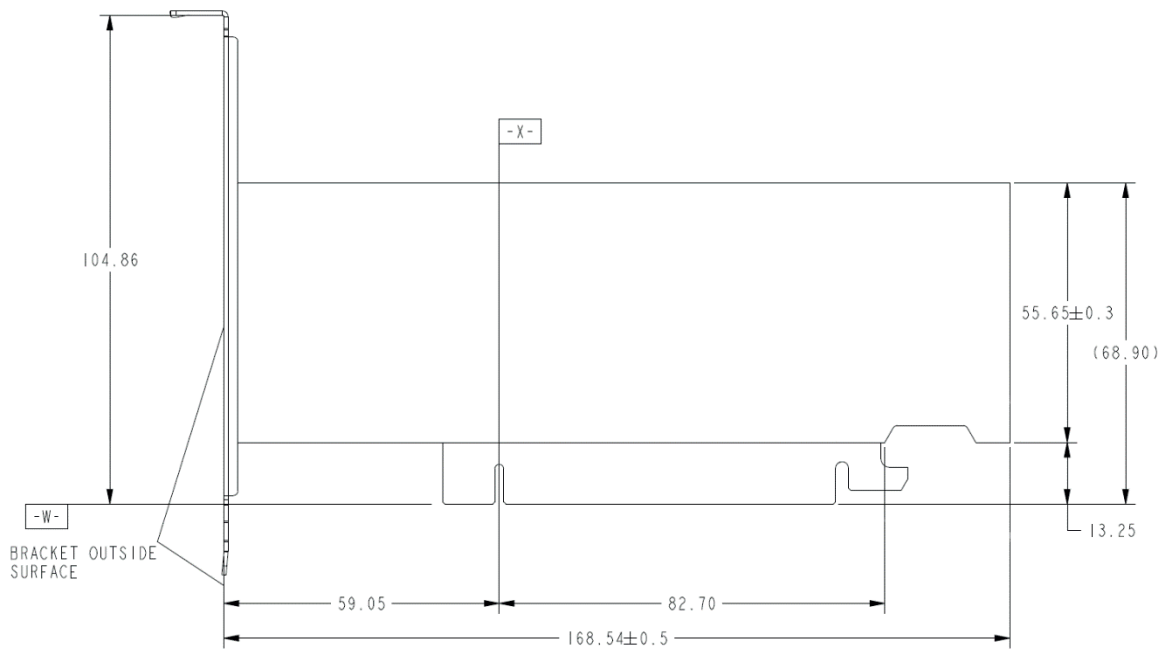
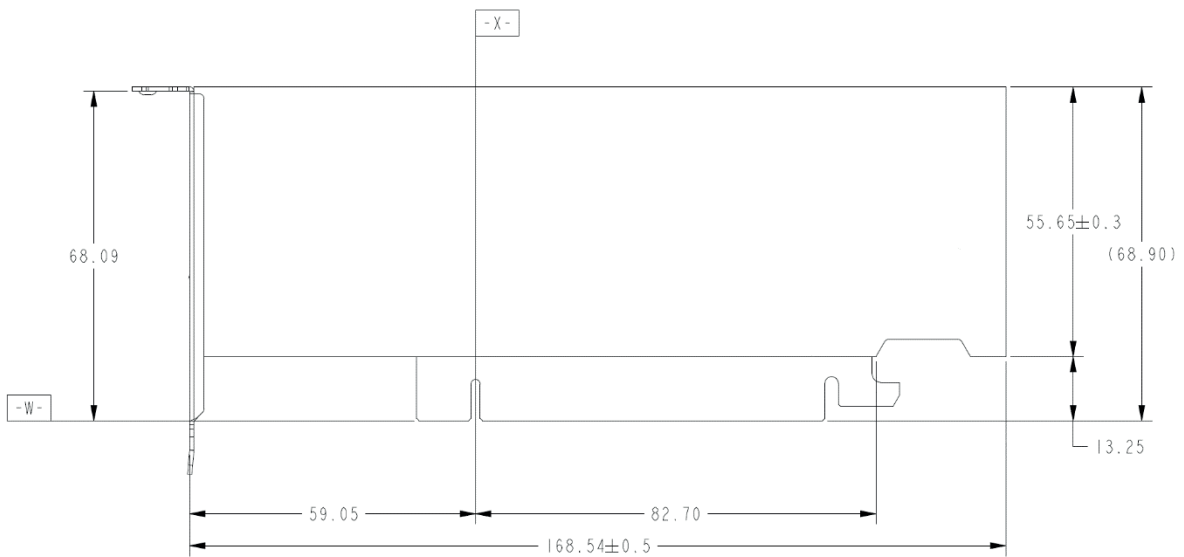


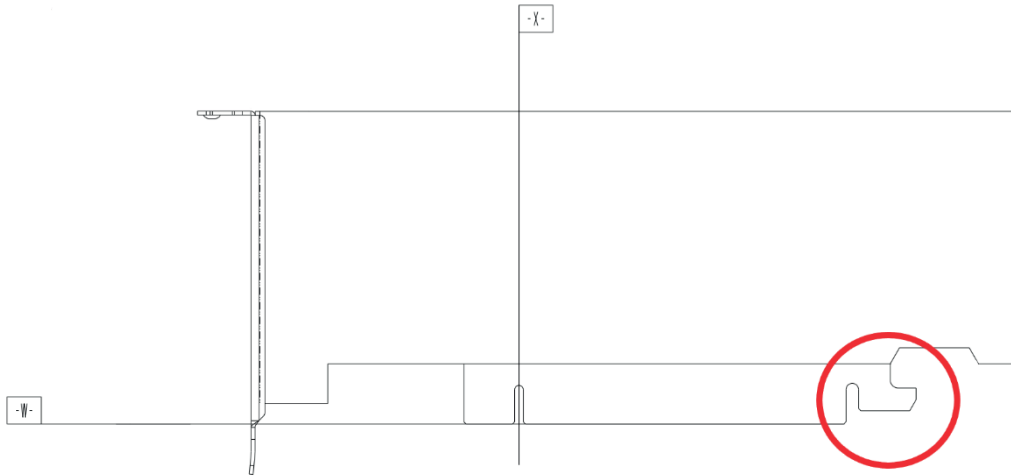
Figure 4. NVIDIA L4 PCIe Card Dimensions with Low Profile Bracket



Hockey Stick Board Retention

The NVIDIA L4 enables south edge board retention using a “hockey stick” tab located to the east of the PCIe card fingers, as shown in Figure 5.

Figure 5. NVIDIA L4 Hockey Stick Tab



Languages

Table 5. Languages Supported

Languages	Windows¹	Linux
English (US)	Yes	Yes
English (UK)	Yes	Yes
Arabic	Yes	
Chinese, Simplified	Yes	
Chinese, Traditional	Yes	
Czech	Yes	
Danish	Yes	
Dutch	Yes	
Finnish	Yes	
French (European)	Yes	
German	Yes	
Greek	Yes	
Hebrew	Yes	
Hungarian	Yes	
Italian	Yes	
Japanese	Yes	
Korean	Yes	
Norwegian	Yes	
Polish	Yes	
Portuguese (Brazil)	Yes	
Portuguese (European/Iberian)	Yes	
Russian	Yes	
Slovak	Yes	
Slovenian	Yes	
Spanish (European)	Yes	
Spanish (Latin America)	Yes	
Swedish	Yes	
Thai	Yes	
Turkish	Yes	
Note:		
¹ Microsoft Windows 10, Windows 11, Windows Server 2019, and Windows Server 2022 Windows are supported.		

Notice

This document is provided for information purposes only and shall not be regarded as a warranty of a certain functionality, condition, or quality of a product. NVIDIA Corporation ("NVIDIA") makes no representations or warranties, expressed or implied, as to the accuracy or completeness of the information contained in this document and assumes no responsibility for any errors contained herein. NVIDIA shall have no liability for the consequences or use of such information or for any infringement of patents or other rights of third parties that may result from its use. This document is not a commitment to develop, release, or deliver any Material (defined below), code, or functionality.

NVIDIA reserves the right to make corrections, modifications, enhancements, improvements, and any other changes to this document, at any time without notice.

Customer should obtain the latest relevant information before placing orders and should verify that such information is current and complete. NVIDIA products are sold subject to the NVIDIA standard terms and conditions of sale supplied at the time of order acknowledgement, unless otherwise agreed in an individual sales agreement signed by authorized representatives of NVIDIA and customer ("Terms of Sale"). NVIDIA hereby expressly objects to applying any customer general terms and conditions with regards to the purchase of the NVIDIA product referenced in this document. No contractual obligations are formed either directly or indirectly by this document.

NVIDIA products are not designed, authorized, or warranted to be suitable for use in medical, military, aircraft, space, or life support equipment, nor in applications where failure or malfunction of the NVIDIA product can reasonably be expected to result in personal injury, death, or property or environmental damage. NVIDIA accepts no liability for inclusion and/or use of NVIDIA products in such equipment or applications and therefore such inclusion and/or use is at customer's own risk.

NVIDIA makes no representation or warranty that products based on this document will be suitable for any specified use. Testing of all parameters of each product is not necessarily performed by NVIDIA. It is customer's sole responsibility to evaluate and determine the applicability of any information contained in this document, ensure the product is suitable and fit for the application planned by customer, and perform 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 document. NVIDIA accepts no 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 document or (ii) customer product designs.

No license, either expressed or implied, is granted under any NVIDIA patent right, copyright, or other NVIDIA intellectual property right under this document. Information published by NVIDIA regarding third-party products or services does not constitute a license from NVIDIA to use such products or services or a warranty or endorsement thereof. Use of such information may require a license from a third party under the patents or other intellectual property rights of the third party, or a license from NVIDIA under the patents or other intellectual property rights of NVIDIA.

Reproduction of information in this document is permissible only if approved in advance by NVIDIA in writing, reproduced without alteration and in full compliance with all applicable export laws and regulations, and accompanied by all associated conditions, limitations, and notices.

THIS DOCUMENT AND 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. TO THE EXTENT NOT PROHIBITED BY LAW, IN NO EVENT WILL NVIDIA BE LIABLE FOR ANY DAMAGES, INCLUDING WITHOUT LIMITATION ANY DIRECT, INDIRECT, SPECIAL, INCIDENTAL, PUNITIVE, OR CONSEQUENTIAL DAMAGES, HOWEVER CAUSED AND REGARDLESS OF THE THEORY OF LIABILITY, ARISING OUT OF ANY USE OF THIS DOCUMENT, EVEN IF NVIDIA HAS BEEN ADVISED OF THE POSSIBILITY OF SUCH DAMAGES. Notwithstanding any damages that customer might incur for any reason whatsoever, NVIDIA's aggregate and cumulative liability towards customer for the products described herein shall be limited in accordance with the Terms of Sale for the product.

Trademarks

NVIDIA, the NVIDIA logo, CUDA, NVIDIA-Certified System, and NVIDIA GPU Boost are trademarks and/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

© 2023 NVIDIA Corporation. All rights reserved.