



NVIDIA RTX A4000

SLEEK DESIGN.
POWERFUL PERFORMANCE.

Amplified Performance for Professionals

The NVIDIA RTX™ A4000 is the most powerful single-slot GPU for professionals, delivering real-time ray tracing, AI-accelerated compute, and high-performance graphics to your desktop. Built on the NVIDIA Ampere architecture, the RTX A4000 combines 48 second-generation RT Cores, 192 third-generation Tensor Cores, and 6,144 CUDA® cores with 16 GB of graphics memory with error-correction code (ECC) so you can innovate with uncompromised computing accuracy and reliability. The RTX A4000 also features a power-efficient, single-slot PCIe form factor that fits into a wide range of workstation chassis, so you can do exceptional work without limits.

NVIDIA RTX professional graphics cards are certified with a broad range of professional applications, tested by leading independent software vendors (ISVs) and workstation manufacturers, and backed by a global team of support specialists. Get the peace of mind needed to focus on what matters with the premier visual computing solution for mission-critical business.

Features

- > PCI Express Gen 4
- > Four DisplayPort 1.4a connectors
- > AV1 decode support
- > DisplayPort with audio
- > 3D stereo support with stereo connector
- > NVIDIA GPUDirect® for Video support
- > NVIDIA Quadro® Sync II¹ compatibility
- > NVIDIA RTX Experience™
- > NVIDIA RTX Desktop Manager software
- > NVIDIA RTX IO support
- > HDCP 2.2 support
- > NVIDIA Mosaic² technology

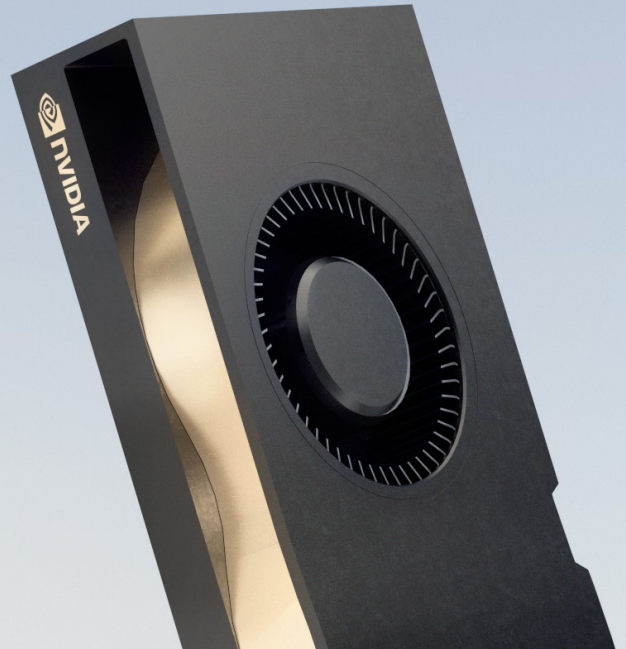
SPECIFICATIONS

Part Number	VCNRTXA4000-PB
EAN Code	3536403383855
GPU memory	16 GB GDDR6
Memory interface	256-bit
Memory bandwidth	448 GB/s
Error-correcting code (ECC)	Yes
NVIDIA Ampere architecture-based CUDA Cores	6,144
NVIDIA third-generation Tensor Cores	192
NVIDIA second-generation RT Cores	48
Single-precision performance	19.2 TFLOPS ³
RT Core performance	37.4 TFLOPS ³
Tensor performance	153.4 TFLOPS ⁴
System interface	PCI Express 4.0 x16
Power consumption	Total board power: 140 W
Thermal solution	Active
Form factor	4.4" H x 9.5" L, single slot
Display connectors	4x DisplayPort 1.4a
Max simultaneous displays	4x 4096 x 2160 @ 120 Hz, 4x 5120 x 2880 @ 60 Hz, 2x 7680 x 4320 @ 60 Hz
Power connector	1x 6-pin PCIe
Encode/decode engines	1x encode, 1x decode (+AV1 decode)
VR ready	Yes
Graphics APIs	DirectX 12.07 ⁵ , Shader Model 5.17 ⁵ , OpenGL 4.68 ⁶ , Vulkan 1.2 ⁶

To learn more about the NVIDIA RTX A4000, visit www.nvidia.com/rtx-a4000/

¹ Quadro Sync II card sold separately. | ² Windows 10 and Linux. | ³ Peak rates based on GPU Boost Clock. | ⁴ Effective teraFLOPS (TFLOPS) using the new sparsity feature. | ⁵ GPU supports DX 12.0 API, hardware feature level 12 + 1. | ⁶ Product is based on a published Khronos specification and is expected to pass the Khronos conformance testing process when available. Current conformance status can be found at www.khronos.org/conformance





NVIDIA RTX A5000

PERFECTLY BALANCED.
BLAZING PERFORMANCE.

Amplified Performance for Professionals

The NVIDIA RTX™ A5000 delivers the power, performance, capabilities, and reliability professionals need to bring their boldest ideas to life. Built on the NVIDIA Ampere architecture, the RTX A5000 combines 64 second-generation RT Cores, 256 third-generation Tensor Cores, and 8,192 CUDA® cores with 24 GB of graphics memory to supercharge rendering, AI, graphics, and compute tasks. Connect two RTX A5000s with NVIDIA NVLink¹ to scale memory and performance with multi-GPU configurations², allowing professionals to work with memory intensive tasks such as large models, ultra-high resolution rendering, and complex compute workloads. Support for NVIDIA virtual GPU software increases the versatility for enterprise deployments.

NVIDIA RTX professional graphics cards are certified with a broad range of professional applications, tested by leading independent software vendors (ISVs) and workstation manufacturers, and backed by a global team of support specialists. Get the peace of mind needed to focus on what matters with the premier visual computing solution for mission-critical business.

Features

- > PCI Express Gen 4
- > Four DisplayPort 1.4a connectors
- > AV1 decode support
- > DisplayPort with audio
- > 3D stereo support with stereo connector
- > NVIDIA GPUDirect® for Video support
- > NVIDIA virtual GPU (vGPU) software support
- > NVIDIA Quadro® Sync II³ compatibility
- > NVIDIA RTX Experience™
- > NVIDIA RTX Desktop Manager software
- > NVIDIA RTX IO support
- > HDCP 2.2 support
- > NVIDIA Mosaic⁴ technology

To learn more about the NVIDIA RTX A5000, visit www.nvidia.com/rtx-a5000/

¹ NVIDIA NVLink sold separately. | ² Connecting two RTX A5000 cards with NVLink to scale performance and memory capacity to 48GB is only possible if your application supports NVLink technology. Please contact your application provider to confirm their support for NVLink. | ³ Quadro Sync II card sold separately. | ⁴ Windows 10 and Linux. | ⁵ Peak rates based on GPU Boost Clock. | ⁶ Effective teraFLOPS (TFLOPS) using the new sparsity feature. | ⁷ Display ports are on by default for RTX A5000. Display ports are not active when using vGPU software. | ⁸ 8 GPU supports DX 12.0 API, hardware feature level 12 + 1. | ⁹ Product is based on a published Khronos specification and is expected to pass the Khronos conformance testing process when available. Current conformance status can be found at www.khronos.org/conformance

SPECIFICATIONS

Part Number	VGNRTXA5000-PB
EAN Code	3536403383817
GPU memory	24 GB GDDR6
Memory interface	384-bit
Memory bandwidth	768 GB/s
Error-correcting code (ECC)	Yes
NVIDIA Ampere architecture-based CUDA Cores	8,192
NVIDIA third-generation Tensor Cores	256
NVIDIA second-generation RT Cores	64
Single-precision performance	27.8 TFLOPS ⁵
RT Core performance	54.2 TFLOPS ⁵
Tensor performance	222.2 TFLOPS ⁶
NVIDIA NVLink	Low profile bridges connect two NVIDIA RTX A5000 GPUs ¹
NVIDIA NVLink bandwidth	112.5 GB/s (bidirectional)
System interface	PCI Express 4.0 x16
Power consumption	Total board power: 230 W
Thermal solution	Active
Form factor	4.4" H x 10.5" L, dual slot, full height
Display connectors	4x DisplayPort 1.4a ⁷
Max simultaneous displays	4x 4096 x 2160 @ 120 Hz, 4x 5120 x 2880 @ 60 Hz, 2x 7680 x 4320 @ 60 Hz
Power connector	1x 8-pin PCIe
Encode/decode engines	1x encode, 2x decode (+AV1 decode)
VR ready	Yes
vGPU software support ⁷	NVIDIA vPC/vApps, NVIDIA RTX Virtual Workstation, NVIDIA Virtual Compute Server
vGPU profiles supported	See the Virtual GPU Licensing Guide
Graphics APIs	DirectX 12.0 ⁸ , Shader Model 5.17 ⁹ , OpenGL 4.6 ⁸ , Vulkan 1.2 ⁹





NVIDIA RTX A6000

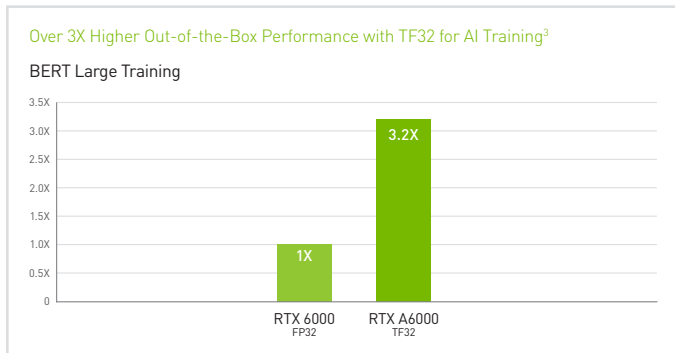
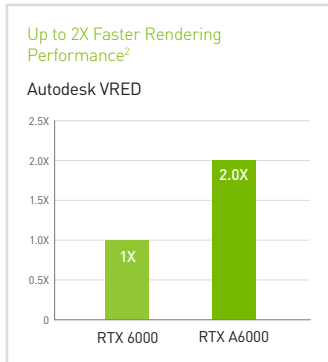
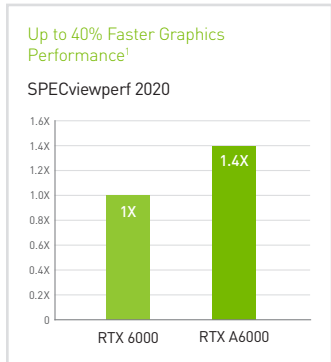
POWERING THE WORLD'S HIGHEST-PERFORMING WORKSTATIONS

Amplified Performance for Professionals

The NVIDIA RTX™ A6000, built on the NVIDIA Ampere architecture, delivers everything designers, engineers, scientists, and artists need to meet the most graphics and compute-intensive workflows. The RTX A6000 is equipped with the latest generation RT Cores, Tensor Cores, and CUDA® cores for unprecedented rendering, AI, graphics, and compute performance. Certified with a broad range of professional applications, tested by leading independent software vendors (ISVs) and workstation manufacturers, and backed by a global team of support specialists, NVIDIA RTX is the visual computing solution of choice for demanding enterprise deployments.

SPECIFICATIONS

GPU memory	48 GB GDDR6
Memory interface	384-bit
Memory bandwidth	768 GB/s
Error-correcting code (ECC)	Yes
NVIDIA Ampere architecture-based CUDA Cores	10,752
NVIDIA third-generation Tensor Cores	336
NVIDIA second-generation RT Cores	84
Single-precision performance	38.7 TFLOPS ⁷
RT Core performance	75.6 TFLOPS ⁷
Tensor performance	309.7 TFLOPS ⁸
NVIDIA NVLink	Connects two NVIDIA RTX A6000 GPUs ¹²
NVIDIA NVLink bandwidth	112.5 GB/s (bidirectional)
System interface	PCI Express 4.0 x16
Power consumption	Total board power: 300 W
Thermal solution	Active
Form factor	4.4" H x 10.5" L, dual slot, full height
Display connectors	4x DisplayPort 1.4a ⁹
Max simultaneous displays	4x 4096 x 2160 @ 120 Hz, 4x 5120 x 2880 @ 60 Hz, 2x 7680 x 4320 @ 60 Hz
Power connector	1x 8-pin CPU
Encode/decode engines	1x encode, 2x decode (+AV1 decode)
VR ready	Yes
vGPU software support	NVIDIA vPC/vApps, NVIDIA RTX Virtual Workstation, NVIDIA Virtual Compute Server
vGPU profiles supported	1 GB, 2 GB, 3 GB, 4 GB, 6 GB, 8 GB, 12 GB, 16 GB, 24 GB, 48 GB
Graphics APIs	DirectX 12.0 ¹⁰ , Shader Model 5.1 ¹⁰ , OpenGL 4.6 ¹¹ , Vulkan 1.18 ¹¹
Compute APIs	CUDA, DirectCompute, OpenCL™

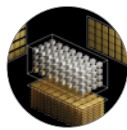


Groundbreaking Innovations



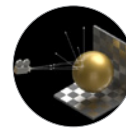
NVIDIA AMPERE ARCHITECTURE

NVIDIA® RTX™ technology revolutionized professional visual computing forever. The NVIDIA Ampere architecture builds on the power of RTX to significantly enhance the performance of rendering, graphics, AI, and compute workloads. Engineered to perfection and featuring cutting-edge innovations, the NVIDIA Ampere architecture takes RTX to new heights for professional workloads.



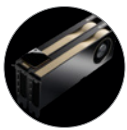
THIRD-GENERATION TENSOR CORES

New Tensor Float 32 (TF32) precision provides up to 5X the training throughput over the previous generation to accelerate AI and data science model training without requiring any code changes. Hardware support for structural sparsity doubles the throughput for inferencing. Tensor Cores also bring AI to graphics with capabilities like DLSS, AI denoising, and enhanced editing for select applications.



SECOND-GENERATION RT CORES

With up to 2X the throughput over the previous generation and the ability to concurrently run ray tracing with either shading or denoising capabilities, second-generation RT Cores deliver massive speedups for workloads like photorealistic rendering of movie content and virtual prototyping of product designs. This technology also speeds up the rendering of ray-traced motion blur for faster results with greater visual accuracy.



THIRD-GENERATION NVLINK

Third-generation NVIDIA NVLink® technology enables users to connect two GPUs together to share GPU performance and memory. With up to 112 gigabytes per second (GB/s) of bidirectional bandwidth and combined graphics memory of up to 96 GB, professionals can tackle the largest rendering, AI, virtual reality, and visual computing workloads. The new NVLink connector also features a shorter Z height, which enables NVLink functionality in a wider range of chassis.



NVIDIA AMPERE ARCHITECTURE-BASED CUDA CORES

The NVIDIA Ampere architecture's CUDA® cores bring double-speed processing for single-precision floating point (FP32) operations and are up to 2X more power efficient than Turing GPUs. This provides significant performance gains for graphics workflows like 3D model development and compute workflows like desktop simulation for computer-aided engineering (CAE).



PCI EXPRESS GEN 4.0

NVIDIA Ampere architecture-based GPUs support PCI Express Gen 4.0 (PCIe Gen 4.0), which provides 2X the bandwidth of PCIe Gen 3.0. This improves data transfer speeds from CPU memory for data-intensive tasks such as AI and data science. Faster PCIe performance also accelerates GPU direct memory access (DMA) transfers, enabling faster video data transfers from GPUDirect® for video-enabled devices and faster input/output (I/O) with GPUDirect Storage.

Features

- > PCI Express Gen 4
- > Four DisplayPort 1.4a connectors
- > AV1 decode support
- > DisplayPort with audio
- > VGA support⁴
- > 3D stereo support with stereo connector
- > NVIDIA GPUDirect® for Video support
- > NVIDIA virtual GPU (vGPU) software support
- > NVIDIA Quadro® Sync II⁵ compatibility
- > NVIDIA Quadro Experience™
- > Desktop Management Software
- > NVIDIA RTX IO support
- > HDCP 2.2 support
- > NVIDIA Mosaic⁶ technology

To learn more about the NVIDIA RTX A6000, visit www.nvidia.com/rtx-a6000

¹ Tests run on workstation with 1x Xeon Gold 6154, 3GHz [3.7GHz Turbo], Win10 x 64, NVIDIA driver version 460.48. SPECviewperf 2020, energy substest. | ² Tests run on workstation with 2x Xeon Gold 6126, 2.6GHz [3.7GHz Turbo], Win10 x 64, NVIDIA driver version 456.37. Autodesk VRED 221.0 GA Release. | ³ Tests run on workstation with AMD Ryzen 3900X, 3.8GHz, 4.6 Turbo, NVIDIA driver 460.17, BERT pre-training throughput using Pytorch, phase 1 sequence length 128, RTX 6000 using FP32 precision, RTX A6000 using TF32 precision. | ⁴ Via adapter/connector/bracket. | ⁵ Quadro Sync II card sold separately. | ⁶ Windows 7, 8, 8.1, 10, and Linux. | ⁷ Peak rates based on GPU Boost Clock. | ⁸ Effective teraFLOPS (TFLOPS) using the new sparsity feature. | ⁹ Display ports are on by default for Quadro RTX A6000. Display ports are not active when using vGPU software. | ¹⁰ GPU supports DX 12.0 API, hardware feature level 12 + 1. | ¹¹ Product is based on a published Khronos specification and is expected to pass the Khronos conformance testing process when available. Current conformance status can be found at www.khronos.org/conformance. | ¹² NVIDIA NVLink sold separately.

