

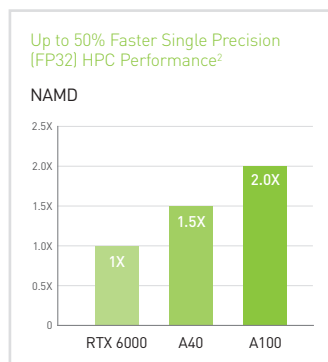
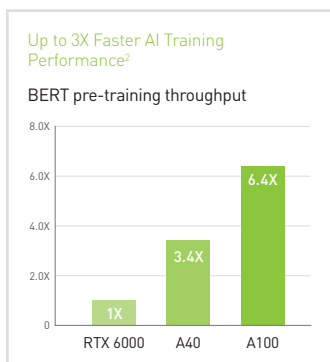
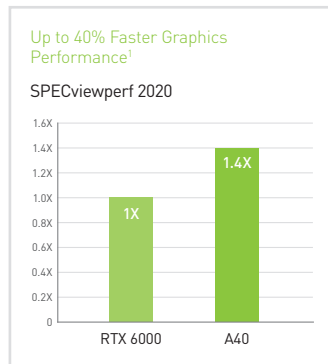
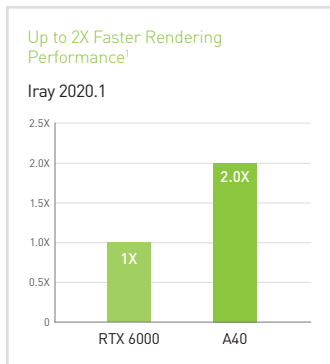
# NVIDIA A40

## Powerful Data Center GPU For Visual Computing

The NVIDIA A40 accelerates the most demanding visual computing workloads from the data center, combining the latest NVIDIA Ampere architecture RT Cores, Tensor Cores, and CUDA® Cores with 48 GB of graphics memory. From powerful virtual workstations accessible from anywhere to dedicated render nodes, NVIDIA A40 brings next-generation NVIDIA RTX™ technology to the data center for the most advanced professional visualization workloads.

### SPECIFICATIONS

GPU architecture	<b>NVIDIA Ampere architecture</b>
GPU memory	<b>48 GB GDDR6 with ECC</b>
Memory bandwidth	<b>696 GB/s</b>
Interconnect interface	<b>NVIDIA® NVLink® 112.5 GB/s [bidirectional]³ PCIe Gen4: 64GB/s</b>
NVIDIA Ampere architecture-based CUDA Cores	<b>10,752</b>
NVIDIA second-generation RT Cores	<b>84</b>
NVIDIA third-generation Tensor Cores	<b>336</b>
Peak FP32 TFLOPS (non-Tensor)	<b>37.4</b>
Peak FP16 Tensor TFLOPS with FP16 Accumulate	<b>149.7   299.4*</b>
Peak TF32 Tensor TFLOPS	<b>74.8   149.6*</b>
RT Core performance TFLOPS	<b>73.1</b>
Peak BF16 Tensor TFLOPS with FP32 Accumulate	<b>149.7   299.4*</b>
Peak INT8 Tensor TOPS	<b>299.3   598.6*</b>
Peak INT4 Tensor TOPS	<b>598.7   1,197.4*</b>
Form factor	<b>4.4" (H) x 10.5" (L) dual slot</b>
Display ports	<b>3x DisplayPort 1.4**; Supports NVIDIA Mosaic and Quadro® Sync⁴</b>
Max power consumption	<b>300 W</b>
Power connector	<b>8-pin CPU</b>
Thermal solution	<b>Passive</b>
Virtual GPU (vGPU) software support	<b>NVIDIA vPC/vApps, NVIDIA RTX Virtual Workstation, NVIDIA Virtual Compute Server</b>
vGPU profiles supported	<b>See the <a href="#">Virtual GPU Licensing Guide</a></b>
NVENC   NVDEC	<b>1x   2x (includes AV1 decode)</b>
Secure and measured boot with hardware root of trust	<b>Yes (optional)</b>
NEBS ready	<b>Level 3</b>
Compute APIs	<b>CUDA, DirectCompute, OpenCL™, OpenACC®</b>
Graphics APIs	<b>DirectX 12.07⁵, Shader Model 5.17⁵, OpenGL 4.68⁶, Vulkan 1.18⁶</b>
MIG support	<b>No</b>



\* Structural sparsity enabled

\*\* A40 is configured for virtualization by default with physical display connectors disabled. The display outputs can be enabled via management software tools.

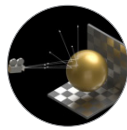
## A Look Inside the NVIDIA Ampere Architecture



### NVIDIA AMPERE ARCHITECTURE CUDA CORES

Double-speed processing for single-precision floating

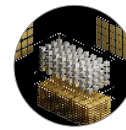
point (FP32) operations and improved power efficiency provide significant performance gains in graphics and compute workflows such as complex 3D computer-aided design (CAD) and computer-aided engineering (CAE).



### 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, architectural design evaluations, 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 TENSOR CORES

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 any code changes. Hardware support for structural sparsity provides up to double the throughput for inferencing. Tensor Cores also bring AI to graphics with capabilities like deep learning super sampling (DLSS), AI denoising, and enhanced editing for select applications.



### 48 GB GDDR6 MEMORY WITH NVLINK

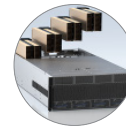
Ultra-fast GDDR6 memory, scalable up to 96 GB with NVLink<sup>3</sup>, gives data scientists,

engineers, and creative professionals the large memory necessary to work with massive datasets and workloads like data science and simulation.



### PCI EXPRESS GEN 4

PCI Express Gen 4 doubles the bandwidth of PCIe Gen 3, improving data-transfer speeds from CPU memory for data-intensive tasks like AI, data science, and 3D design. Faster PCIe performance also accelerates GPU direct memory access (DMA) transfers, providing faster input/output communication of video data between the GPU and GPUdirect<sup>®</sup> for Video-enabled devices to deliver a powerful solution for live broadcast. A40 is backwards compatible with PCI Express Gen 3 for deployment flexibility.



### DATA CENTER EFFICIENCY AND SECURITY

Featuring a dual-slot, power-efficient design, NVIDIA A40 is up to 2X as power efficient as the previous generation and compatible with a wide range of servers from worldwide OEMs. The NVIDIA A40 includes secure and measured boot with hardware root-of-trust technology, ensuring that firmware isn't tampered with or corrupted.

The NVIDIA A40 GPU delivers state-of-the-art visual computing capabilities, including real-time ray tracing, AI acceleration, and multi-workload flexibility to accelerate deep learning, data science, and compute-based workloads. Virtual workstations powered by NVIDIA A40 and NVIDIA RTX Virtual Workstation (vWS) and NVIDIA Virtual Compute Server software benefit from extensive testing across a broad range of industry applications and professional software for optimal performance and stability.

#### EVERY DEEP LEARNING FRAMEWORK

mxnet

PYTORCH

APACHE SPARK

TensorFlow

#### RTX FOR PROFESSIONAL APPLICATIONS

Pr Adobe Premiere Pro

SOLIDWORKS

PLM Software  
SIEMENS  
NX

AUTODESK  
ARNOLD



REDSHIFT

AUTODESK  
VRED

KeyShot

UNREAL  
ENGINE

blender

octane-render

v-ray

### Learn more

To learn more about the NVIDIA A40 GPU, visit [www.nvidia.com/a40](http://www.nvidia.com/a40)

1 Rendering and Graphics tests run on 2x Xeon Gold 6126 2.6GHz (3.7GHz Turbo). 256GB system memory. NVIDIA Driver 461.09. Rendering test: Iray 2020.1. Render time of NVIDIA Endeavor scene. Graphics test: SPECviewperf 2020 Subtest, 4K medical-03 Composite. | 2 AI and HPC tests run on AMD EPYC 7742@2.25GHz (3.4GHz Turbo). 512GB system memory. NVIDIA Driver 460.14. AI Training: BERT pre-training throughput. PyTorch (2/3) Phase 1 and (1/3) Phase 2. Precision FP32 for RTX 6000 and TF32 for A40 and A100. Sequence length for Phase 1 = 128. Phase 2 = 512. Single Precision HPC: NAMD version 3.0a7, stmv\_nve\_cuda; Precision=FP32; ns/day, CUDA Version: 11.1.74. | 3 Connecting two NVIDIA A40 cards with NVLink to scale performance and memory capacity to 96 GB is only possible if your application supports NVLink technology. Please contact your application provider to confirm their support for NVLink. | 4 Quadro Sync II card sold separately. Mosaic supported on Windows 10 and Linux. | 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](http://www.khronos.org/conformance)

