



NVIDIA TESLA ONE PLATFORM. UNLIMITED DATA CENTER ACCELERATION.






The Exponential Growth of Computing

Accelerating scientific discovery, visualizing big data for insights, and providing smart AI-based services to consumers are everyday challenges for researchers and engineers. Solving these challenges takes increasingly complex and precise simulations, the processing of tremendous amounts of data, or training and running sophisticated deep learning networks. These workloads also require accelerating data centers to meet the growing demand for exponential computing.

NVIDIA® Tesla® is the world's leading platform for the accelerated data center, deployed by the largest supercomputing centers and enterprises. It enables breakthrough performance with fewer, more powerful servers, resulting in faster scientific discoveries and insights while saving money.

With over 550 high performance computing (HPC) applications GPU-optimized in a broad range of domains, including 15 of the top 15 HPC applications, and all deep learning frameworks, every modern data center can save money with the Tesla platform.

Choose the Right NVIDIA Data Center Product for You.

<p>NVIDIA Tesla V100 with NVIDIA NVLink™</p> 	<p>NVIDIA Tesla V100 PCIe</p> 	<p>NVIDIA Tesla P4</p> 	<p>NVIDIA Tesla P40</p> 	<p>NVIDIA Tesla P6</p> 
<p>DESIGNED FOR Deep Learning</p>	<p>DESIGNED FOR HPC and Deep Learning</p>	<p>DESIGNED FOR Deep Learning Inference and Video Transcoding</p>	<p>DESIGNED FOR GPU Virtualization - Graphics and Compute</p>	<p>DESIGNED FOR GPU Virtualization - Graphics and Compute</p>
<p>Up to 3X faster time to solution over P100</p>	<p>Up to 5X lower total cost of ownership (TCO) than CPUs for mixed workloads</p>	<p>40X higher energy efficiency than CPUs for inference</p>	<p>Up to 24 virtual GPUs per board</p>	<p>Up to 16 virtual GPUs per board</p>
<p>Ultimate deep learning training performance</p>	<p>Most versatility for mixed HPC workloads 32 GB memory configuration for memory-intensive HPC applications</p>	<p>Low power, low profile optimized for scale-out deep learning inference deployment</p>	<p>Industry's highest graphics performance for virtualized environments Run multiple virtualized graphics and compute workloads</p>	<p>Maximum performance for any virtualized workload in a blade-optimized form factor Double the frame buffer of previous-generation NVIDIA Maxwell™</p>
<p>KEY FEATURES</p> <ul style="list-style-type: none"> > 125 teraFLOPS of tensor operations for deep learning > 15.7 teraFLOPS of single-precision performance > 7.8 teraFLOPS of double-precision performance > 300 GB/s NVIDIA NVLink interconnect > 900 GB/s memory bandwidth > 32 GB / 16 GB HBM2 memory 	<p>KEY FEATURES</p> <ul style="list-style-type: none"> > 112 teraFLOPS of tensor operations for deep learning > 14 teraFLOPS of single-precision performance > 7 teraFLOPS of double-precision performance > 900 GB/s memory bandwidth > 32 GB / 16 GB HBM2 memory 	<p>KEY FEATURES</p> <ul style="list-style-type: none"> > 22 teraFLOPS of INT8 inference performance > 5.5 teraFLOPS of single-precision performance > 1 decode and 2 encode video engines > 50 W / 75 W power > Low-profile form factor 	<p>KEY FEATURES</p> <ul style="list-style-type: none"> > 24 GB memory > 24 H.264 1080p30 streams > Up to 24 vGPU instances > PCIe 3.0 dual-slot form factor > 250 W power 	<p>KEY FEATURES</p> <ul style="list-style-type: none"> > 16 GB memory > 24 H.264 1080p30 streams > Up to 16 vGPU instances > MXM form factor > 90 W (70 W opt) power
<p>RECOMMENDED SERVER CONFIGURATIONS</p> <p>8-way NVIDIA NVLink hybrid cube mesh (NVIDIA HGX)</p>	<p>RECOMMENDED SERVER CONFIGURATIONS</p> <p>2-4 GPUs per node</p>	<p>RECOMMENDED SERVER CONFIGURATIONS</p> <p>1-2 GPUs per node</p>	<p>RECOMMENDED SERVER CONFIGURATIONS</p> <p>2-4 GPUs per node</p>	<p>RECOMMENDED SERVER CONFIGURATIONS</p> <p>GPUs per node dependent on the blade server</p>