



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.

Inspired by the demands of deep learning and analytics, NVIDIA® DGX™ Systems are built on the NVIDIA Tesla V100 platform. Combined with innovative GPU-optimized software and simplified management, these fully integrated solutions deliver groundbreaking performance and results.

With over 500 HPC applications GPU-optimized in a broad range of domains, including 10 of the top 10 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

NVIDIA Tesla V100 with NVIDIA NVLink	NVIDIA Tesla V100 PCIe	NVIDIA Tesla P4	NVIDIA Tesla P40	NVIDIA Tesla P6	NVIDIA DGX Station™	NVIDIA DGX-1™
 <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>DESIGNED FOR Full Integrated Deep Learning Workstation</p>	 <p>DESIGNED FOR Full Integrated Deep Learning Server</p>
<p>Up to 3X faster time-to-solution over P100</p>	<p>Up to 4X higher throughput 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>30% faster training with fully-integrated NVIDIA DGX Software</p>	
<p>Ultimate deep learning training performance</p>	<p>Most versatility for mixed HPC workloads</p>	<p>Low power, low profile optimized for scale out deep learning inference deployment</p>	<p>Industry's highest graphics performance for virtualized environments</p> <p>Run multiple virtualized graphics and compute workloads</p>	<p>Maximum performance for any virtualized workload in a blade optimized from factor</p> <p>Double the frame buffer of previous generation NVIDIA Maxwell™</p>	<p>World's first 4-way NVIDIA NVLink Workstation powered by NVIDIA Tesla V100</p>	<p>1 PetaFLOPS data center server powered by NVIDIA Tesla V100 with NVLink</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 half-precision performance > 300 GB/s NVIDIA NVLink™ Interconnect > 900 GB/s memory bandwidth > 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 half-precision performance > 900 GB/s memory bandwidth > 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>KEY FEATURES</p> <ul style="list-style-type: none"> > Workstation form factor > NVIDIA DGX software stack > 500 TeraFLOPS of tensor operations for deep learning > 64 GB system GPU memory > 200 GB/s NVIDIA NVLink bandwidth > 1,500 W power 	<p>KEY FEATURES</p> <ul style="list-style-type: none"> > 3RU server form factor > NVIDIA DGX software stack > 1 PetaFLOPS of tensor operations for deep learning > 128 GB system GPU memory > 300 GB/s NVIDIA NVLink bandwidth > 3,200 W power
<p>RECOMMENDED SERVER CONFIGURATIONS</p> <p>8-way NVIDIA NVLink hybrid cube mesh (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>	<p>GPU INTERCONNECT CONFIGURATION</p> <p>4-way NVIDIA NVLink</p>	<p>GPU INTERCONNECT CONFIGURATION</p> <p>8-way NVIDIA NVLink hybrid cube mesh</p>