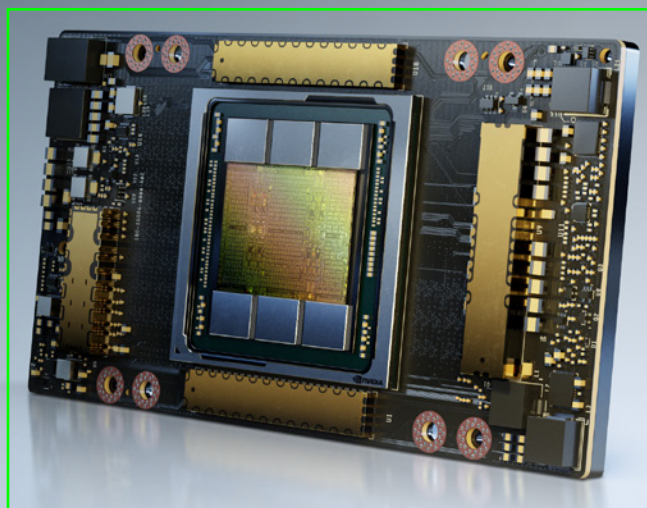




# NVIDIA A100 TENSOR CORE GPU

Unprecedented Acceleration at Every Scale



## The Most Powerful Compute Platform for Every Workload

The NVIDIA A100 Tensor Core GPU delivers unprecedented acceleration—at every scale—to power the world's highest-performing elastic data centers for AI, data analytics, and high-performance computing (HPC) applications. As the engine of the NVIDIA data center platform, A100 provides up to 20X higher performance over the prior NVIDIA Volta™ generation. A100 can efficiently scale up or be partitioned into seven isolated GPU instances with Multi-Instance GPU (MIG), providing a unified platform that enables elastic data centers to dynamically adjust to shifting workload demands.

NVIDIA A100 Tensor Core technology supports a broad range of math precisions, providing a single accelerator for every workload. The latest generation A100 80GB doubles GPU memory and debuts the world's fastest memory bandwidth at 2 terabytes per second (TB/s), speeding time to solution for the largest models and most massive datasets.

A100 is part of the complete NVIDIA data center solution that incorporates building blocks across hardware, networking, software, libraries, and optimized AI models and applications from the NVIDIA NGC™ catalog. Representing the most powerful end-to-end AI and HPC platform for data centers, it allows researchers to deliver real-world results and deploy solutions into production at scale.

## NVIDIA A100 TENSOR CORE GPU SPECIFICATIONS (SXM4 AND PCIe FORM FACTORS)

	A100 40GB PCIe	A100 80GB PCIe	A100 40GB SXM	A100 80GB SXM
FP64	9.7 TFLOPS			
FP64 Tensor Core	19.5 TFLOPS			
FP32	19.5 TFLOPS			
Tensor Float 32 (TF32)	156 TFLOPS   312 TFLOPS*			
BFLOAT16 Tensor Core	312 TFLOPS   624 TFLOPS*			
FP16 Tensor Core	312 TFLOPS   624 TFLOPS*			
INT8 Tensor Core	624 TOPS   1248 TOPS*			
GPU Memory	40GB HBM2	80GB HBM2e	40GB HBM2	80GB HBM2e
GPU Memory Bandwidth	1,555GB/s	1,935GB/s	1,555GB/s	2,039GB/s
Max Thermal Design Power (TDP)	250W	300W	400W	400W
Multi-Instance GPU	Up to 7 MIGs @ 5GB	Up to 7 MIGs @ 10GB	Up to 7 MIGs @ 5GB	Up to 7 MIGs @ 10GB
Form Factor	PCIe		SXM	
Interconnect	NVIDIA® NVLink® Bridge for 2 GPUs: 600GB/s ** PCIe Gen4: 64GB/s		NVLink: 600GB/s PCIe Gen4: 64GB/s	
Server Options	Partner and NVIDIA-Certified Systems™ with 1-8 GPUs		NVIDIA HGX™ A100-Partner and NVIDIA-Certified Systems with 4, 8, or 16 GPUs NVIDIA DGX™ A100 with 8 GPUs	

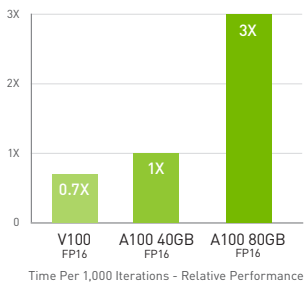
\* With sparsity

\*\* SXM4 GPUs via HGX A100 server boards; PCIe GPUs via NVLink Bridge for up to two GPUs

## Incredible Performance Across Workloads

### Up to 3X Higher AI Training on Largest Models

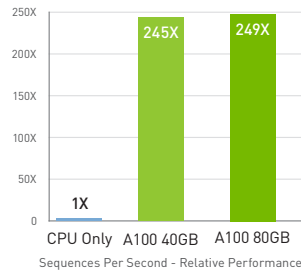
#### DLRM Training



DLRM on HiggsCTR framework, precision = FP16, 1 NVIDIA A100 80GB batch size = 48 | NVIDIA A100 40GB batch size = 32 | NVIDIA V100 32GB batch size = 32

### Up to 249X Higher AI Inference Performance over CPUs

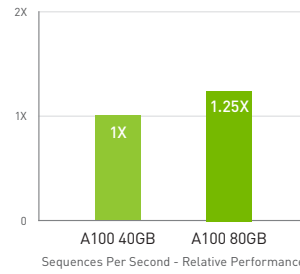
#### BERT-LARGE Inference



BERT-Large Inference | CPU only: Dual Xeon Gold 6240 @2.60 GHz, precision = FP32, batch size = 128 | V100: NVIDIA TensorRT™ (TRT) 7.2, precision = INT8, batch size = 256 | A100 40GB and 80GB, batch size = 256, precision = INT8 with sparsity

### Up to 1.25X Higher AI Inference Performance over A100 40GB

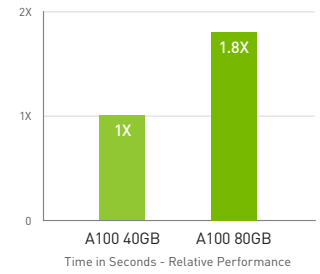
#### RNN-T Inference: Single Stream



MLPerf 0.7 RNN-T measured with 1/7 MIG slices. Framework: TensorRT 7.2, dataset = LibriSpeech, precision = FP16

### Up to 1.8X Higher Performance for HPC Applications

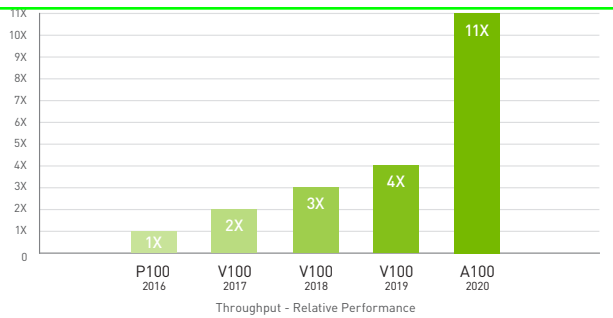
#### Quantum Espresso



Quantum Espresso measured using CNT10P0R8 dataset, precision = FP64

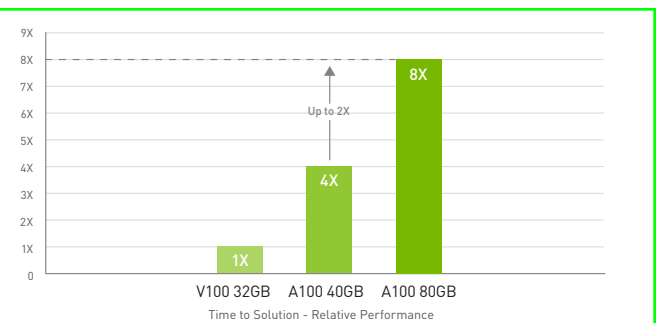
### 11X More HPC Performance in Four Years

#### Throughput for Top HPC Apps



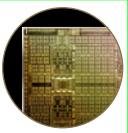
Geometric mean of application speedsups vs. P100. Benchmark application: Amber [PME-Cellulose\_NVE], Chroma [szsc121\_24\_128], GROMACS [ADH Dodec], MILC [Apex Medium], NAMD [stmv\_nve\_cuda], PyTorch [BERT-Large Fine Tuner], Quantum Espresso [AUSURF112-JR], Random Forest FP32 [make\_blobs 160000 x 64: 10], TensorFlow [ResNet-50], VASP & [5i Hugs] | GPU node with dual-socket CPUs with 4x NVIDIA P100, V100, or A100 GPUs

### 2X Faster than A100 40GB on Big Data Analytics Benchmark



Big data analytics benchmark | GPU-BDB is derived from the TPCx-BB benchmark and is used for internal performance testing. Results from GPU-BDB are not comparable to TPCx-BB | 30 analytical retail queries, ETL, ML, NLP on 10TB dataset | V100 32GB, RAPIDS/Dask | A100 40GB and A100 80GB, RAPIDS/Dask/BlazingSQL

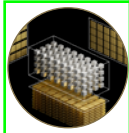
## Groundbreaking Innovations



### NVIDIA AMPERE ARCHITECTURE

Whether using MIG to partition an A100 GPU into smaller instances or NVLink to connect multiple

GPUs to speed large-scale workloads, A100 can readily handle different-sized acceleration needs, from the smallest job to the biggest multi-node workload. A100's versatility means IT managers can maximize the utility of every GPU in their data center, around the clock.



### THIRD-GENERATION TENSOR CORES

NVIDIA A100 delivers 312 teraFLOPS (TFLOPS) of deep learning performance. That's 20X

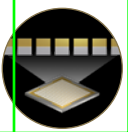
the Tensor floating-point operations per second (FLOPS) for deep learning training and 20X the Tensor tera operations per second (TOPS) for deep learning inference compared to NVIDIA Volta GPUs.



### NEXT-GENERATION NVLINK

NVIDIA NVLink in A100 delivers 2X higher throughput compared to the previous generation. When combined with NVIDIA NVSwitch™,

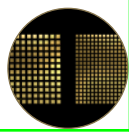
up to 16 A100 GPUs can be interconnected at up to 600 gigabytes per second (GB/sec), unleashing the highest application performance possible on a single server. NVLink is available in A100 SXM GPUs via HGX A100 server boards and in PCIe GPUs via an NVLink Bridge for up to 2 GPUs.



### MULTI-INSTANCE GPU (MIG)

An A100 GPU can be partitioned into as many as seven GPU instances, fully isolated at the hardware level with their

own high-bandwidth memory, cache, and compute cores. MIG gives developers access to breakthrough acceleration for all their applications, and IT administrators can offer right-sized GPU acceleration for every job, optimizing utilization and expanding access to every user and application.



### HIGH-BANDWIDTH MEMORY (HBM2E)

With up to 80 gigabytes of HBM2e, A100 delivers the world's fastest GPU memory bandwidth

of over 2TB/s, as well as a dynamic random-access memory (DRAM) utilization efficiency of 95%. A100 delivers 1.7X higher memory bandwidth over the previous generation.



### STRUCTURAL SPARSITY

AI networks have millions to billions of parameters. Not all of these parameters are needed for accurate predictions, and some

can be converted to zeros, making the models "sparse" without compromising accuracy. Tensor Cores in A100 can provide up to 2X higher performance for sparse models. While the sparsity feature more readily benefits AI inference, it can also improve the performance of model training.

The NVIDIA A100 Tensor Core GPU is the flagship product of the NVIDIA data center platform for deep learning, HPC, and data analytics. The platform accelerates over 2,000 applications, including every major deep learning framework. A100 is available everywhere, from desktops to servers to cloud services, delivering both dramatic performance gains and cost-saving opportunities.

## OPTIMIZED SOFTWARE AND SERVICES FOR ENTERPRISE



### EVERY DEEP LEARNING FRAMEWORK

*mxnet*

PYTORCH



 TensorFlow

### 2,000+ GPU-ACCELERATED APPLICATIONS



Altair nanoFluidX



Altair ultraFluidX



AMBER



ANSYS Fluent



DS SIMULIA Abaqus



GAUSSIAN



GROMACS



NAMD



OpenFOAM



VASP



WRF

To learn more about the NVIDIA A100 Tensor Core GPU, visit [www.nvidia.com/a100](https://www.nvidia.com/a100)

© 2021 NVIDIA Corporation. All rights reserved. NVIDIA, the NVIDIA logo, DGX, HGX, NGC, NVIDIA-Certified Systems, NVLink, NVSwitch, and Volta 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. All other trademarks are property of their respective owners. JUN21

