



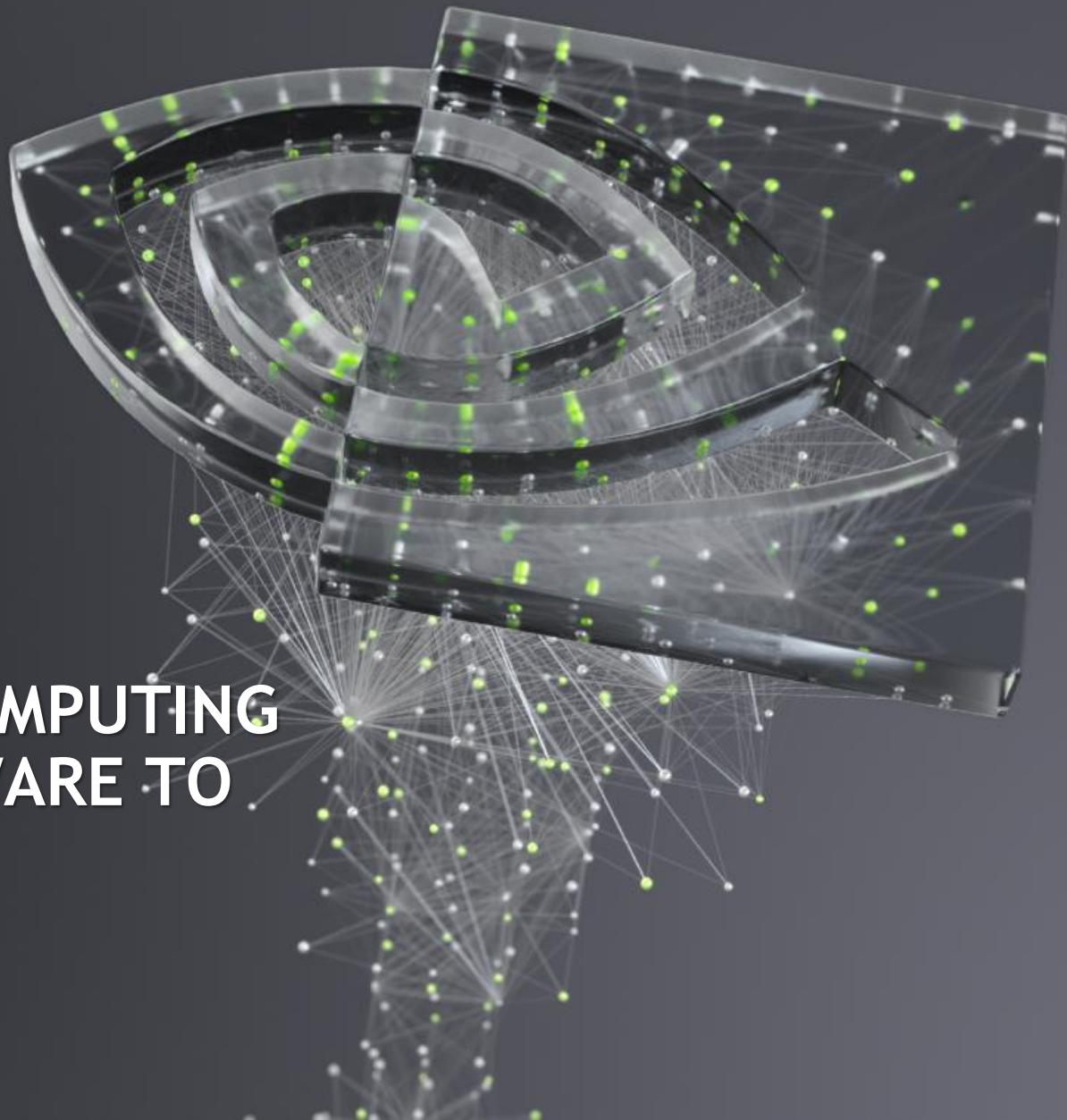
TECH TAKE

NVIDIA INSIGHTS: GPU COMPUTING LANDSCAPE FROM SOFTWARE TO HARDWARE.

Michael Lang

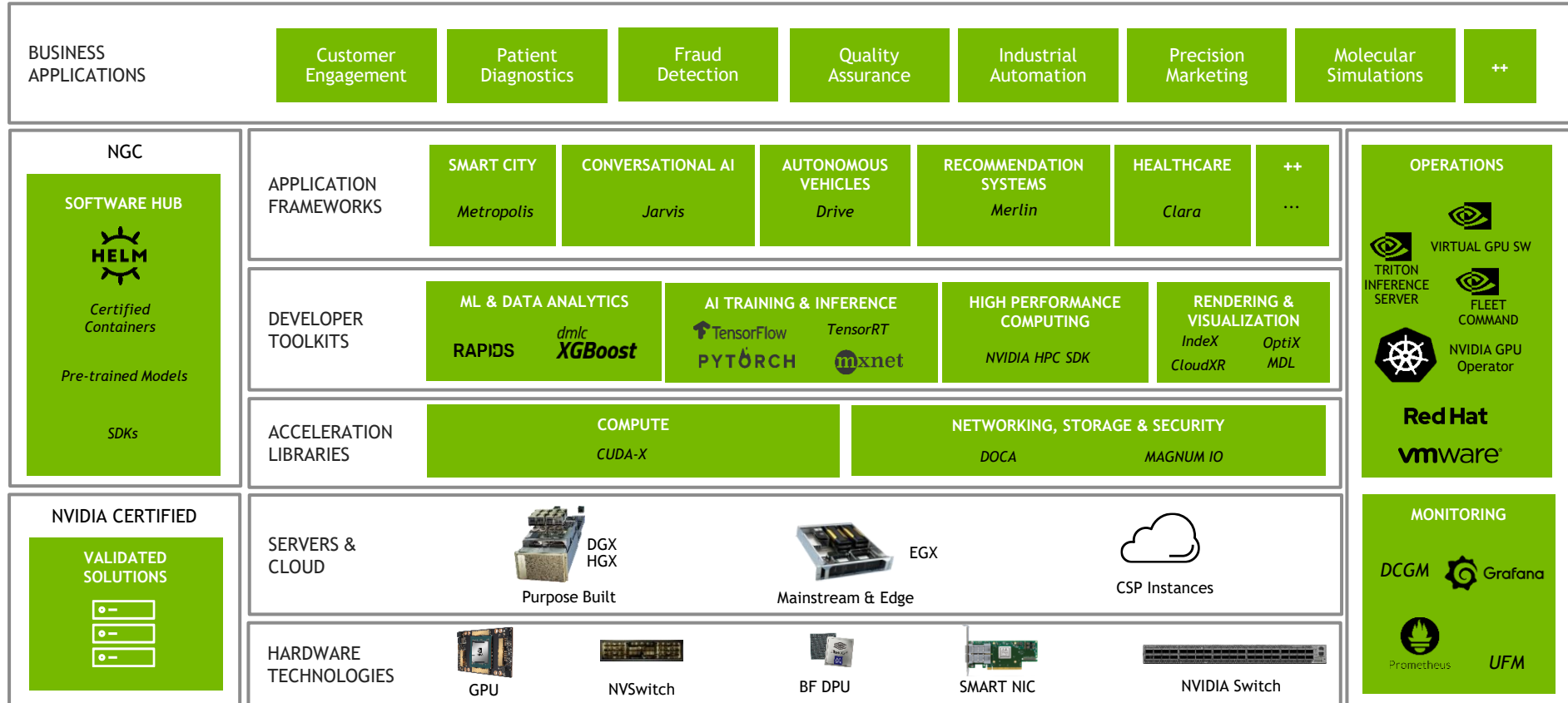
Solutions Architecture Manager - APAC South

June 29 2021

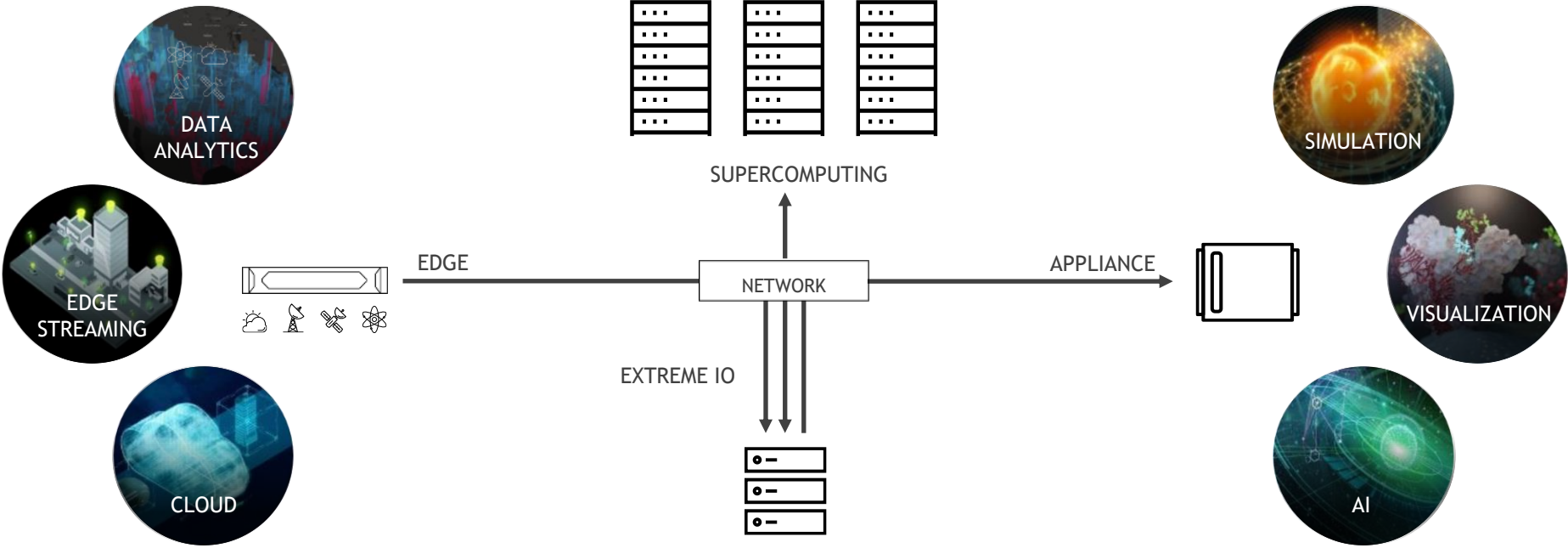


NVIDIA IS A FULL STACK COMPUTING PLATFORM

Amazing Innovation and Expansion of NVIDIA Ecosystem



EXPANDING UNIVERSE OF SCIENTIFIC COMPUTING



NVIDIA POWER WORLD'S FASTEST AND MOST EFFICIENT, AI SUPERCOMPUTERS



8

of Top 10

68%

Overall

20%

More InfiniBand Systems v. ISC20



9

of Top 10

3.5X

Higher Energy Efficiency v.
non-GPU Green500 Systems

29.5GF/W

Greenest NVIDIA System



8

of Top 10

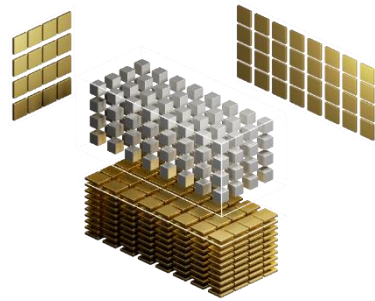
1.1EF

Mixed Precision AI
Performance on SUMMIT

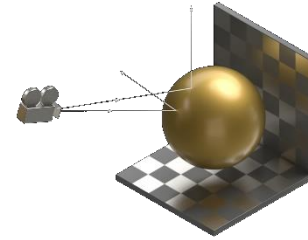
2X

Submission Over Previous List

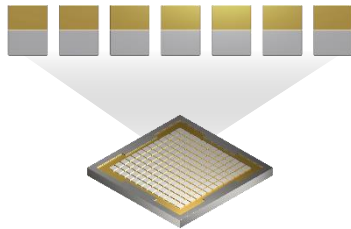
AMPERE ARCHITECTURE



3rd Gen Tensor Cores
Faster, Flexible, Easier to use
20x AI Perf with TF32
2.5x HPC Perf



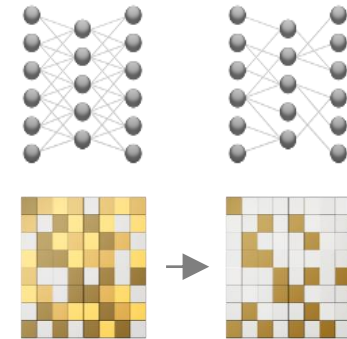
2nd Gen RT Cores
Up to 2X throughput of
previous generation



New Multi-Instance GPU
Optimal utilization with right sized GPU
7x Simultaneous Instances per GPU



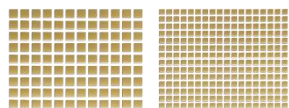
3rd Gen NVLINK and NVSWITCH
Efficient Scaling to Enable Super GPU
2X More Bandwidth



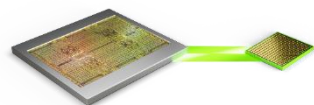
Sparsity
Harness Sparsity in AI Models
2x AI Performance

ANNOUNCING NVIDIA A100 80GB

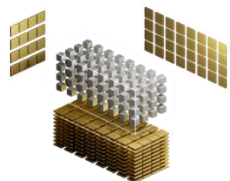
Supercharging The World's Highest Performing AI Supercomputing GPU



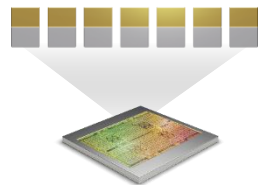
80GB HBM2e
For largest datasets and models



2TB/s +
World's highest memory bandwidth to
feed the world's fastest GPU



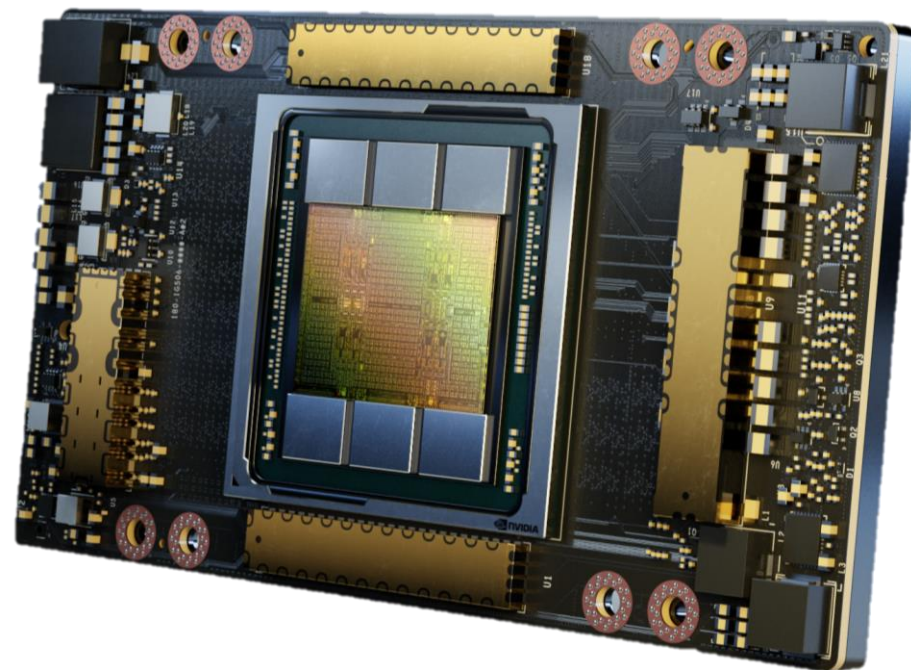
3rd Gen Tensor Core



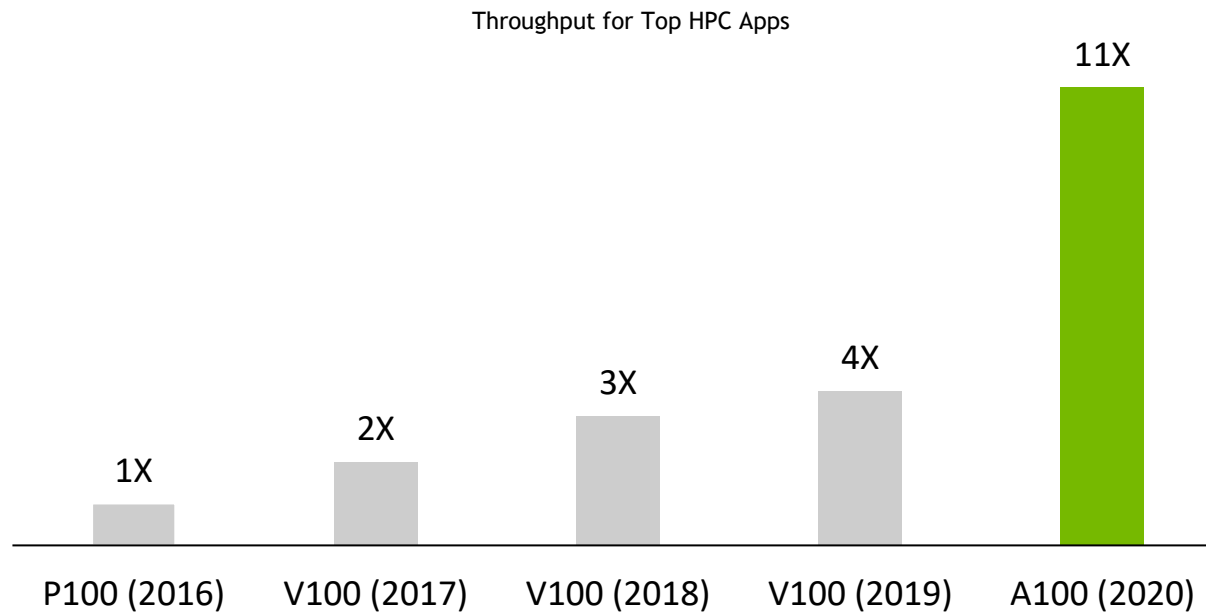
Multi-Instance GPU



3rd Gen NVLink



11X MORE HPC PERFORMANCE IN FOUR YEARS

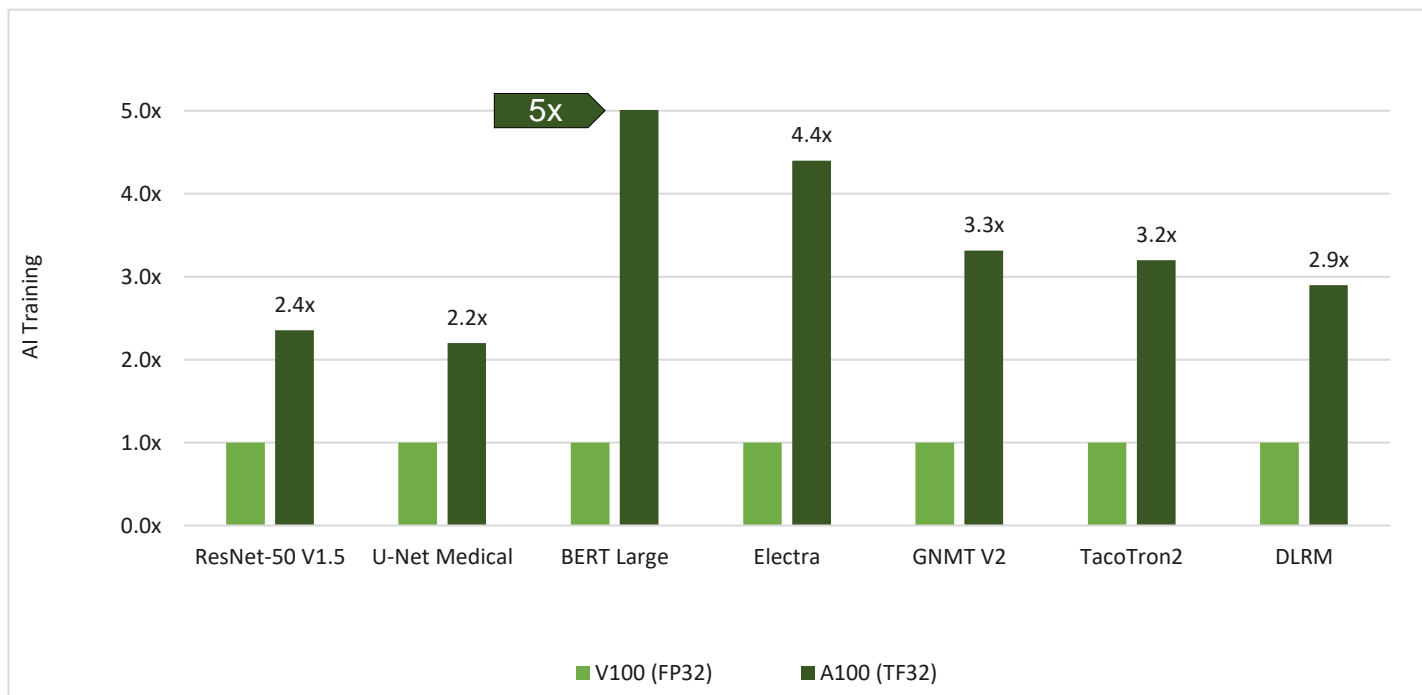


Geometric mean of application speedups vs. P100: Benchmark application: Amber [PME-Cellulose_NVE], Chroma [szscl21_24_128], GROMACS [ADH Dodec], MILC [Apex Medium], NAMD [stmv_nve_cuda], PyTorch (BERT-Large Fine Tuner), Quantum Espresso [AUSURF112-jR]; Random Forest FP32 [make_blobs (160000 x 64 : 10)], TensorFlow [ResNet-50], VASP 6 [Si Huge] | GPU node with dual-socket CPUs with 4x NVIDIA P100, V100, or A100 GPUs.

3RD GENERATION MULTI-PRECISION TENSOR CORES IN A100

Greatest Generational Leap - 20X Volta

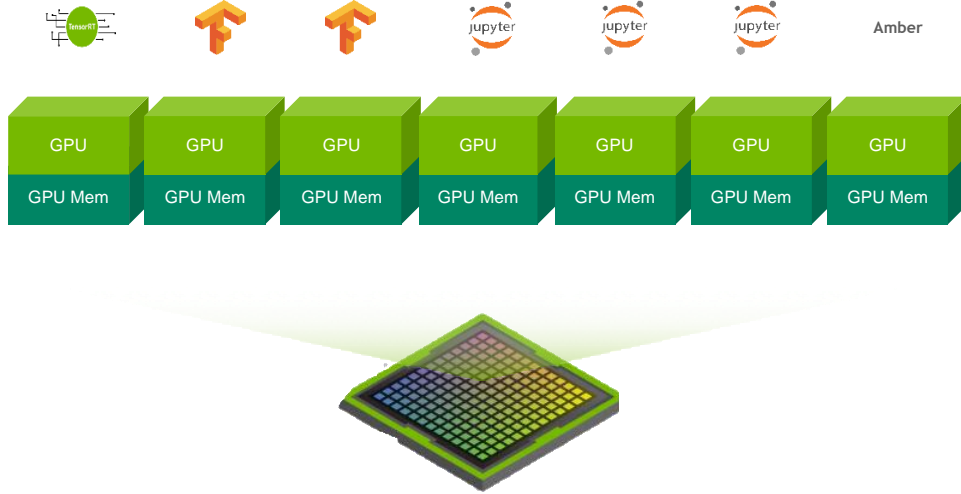
Precision	Peak	Vs Volta
FP64 HPC	19.5 TFLOPS	2.5x
FP32 TRAINING	312 TFLOPS	20x
INT8 INFERENCE	1,248 TOPS	20x



TF32 = Default Precision in TensorFlow, PyTorch, MXNet

NVIDIA A100 MULTI-INSTANCE GPU

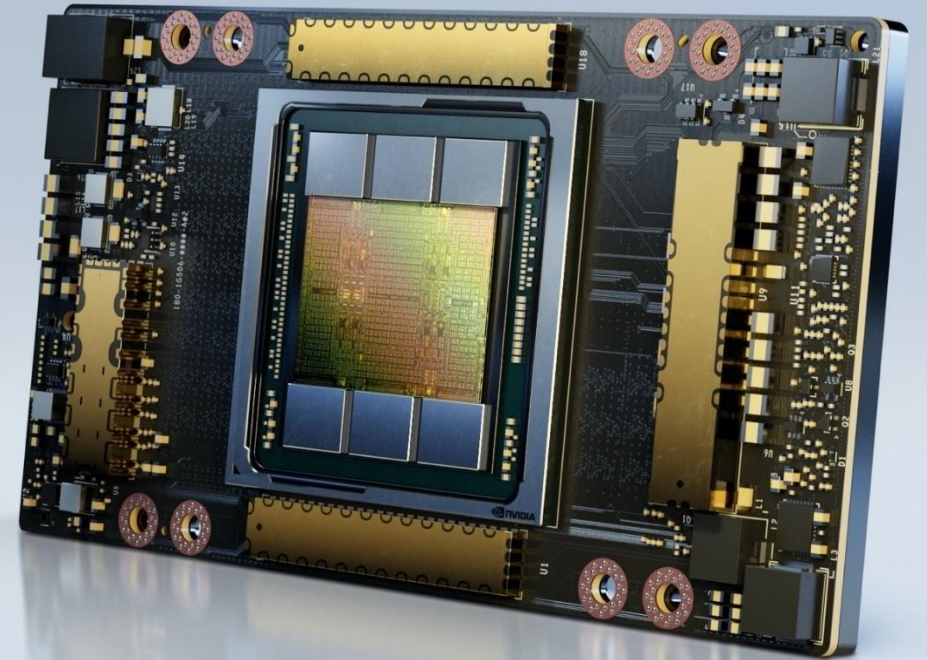
Greatest Generational Leap - 20X Volta



Up To 7 hardware isolated GPU Instances per A100

Simultaneous Workload Execution With Guaranteed Quality Of Service

Right-sized GPU



A100 AVAILABLE VIA NVIDIA HGX A100 AND A100 PCIe

A100 PCIe



For Mainstream Servers

1-8 GPUs per server, optional NVLink Bridge between 2 GPUs

40GB option only

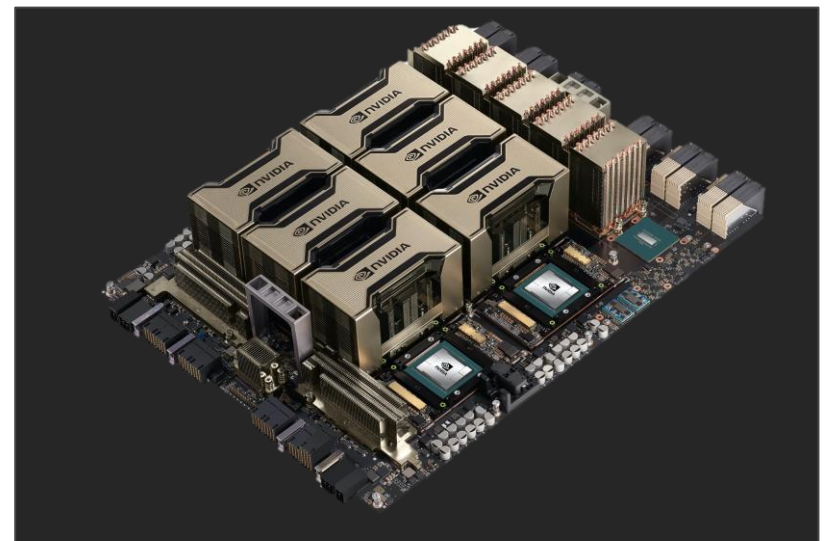
HGX A100 4-GPU



Scale-Up - Mixed AI & HPC

4 A100s, Fully Connected w/ shared NVLinks

HGX A100 8-GPU

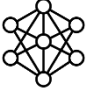
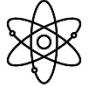




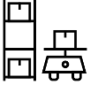






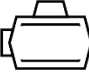





Scale-up - Fastest Time-to-solution for AI

8 GPUs, Full NVLink B/W between all GPUs with NVSwitch

GPUs FOR EVERY VIRTUAL WORKLOAD

Expanding Workloads Drive The Need for Specialized Accelerators

 <p>DL Training</p>  <p>Scientific Research</p>  <p>Data Analytics</p> <hr/> <p>Highest Perf Compute AI, HPC, Data Processing</p> <p>Fastest Compute, FP64 Up to 7 MIG instances</p>	 <p>Language Processing</p>  <p>Conversational AI</p>  <p>Recommender Systems</p> <hr/> <p>AI Inference & Mainstream Compute</p> <p>Versatile Mainstream Compute FP64, Up to 4 MIG instances</p>	 <p>Edge AI</p>  <p>Edge Video</p>  <p>Mobile Cloud Games</p> <hr/> <p>Small Footprint Datacenter and Edge Inference</p> <p>High density Video & Graphics Compact & Versatile</p>	 <p>Virtual Workstation</p>  <p>Video Conferencing</p>  <p>4K Cloud Games</p> <hr/> <p>Mainstream Graphics & Video with AI</p> <p>4K Cloud Gaming Graphics, Video with AI</p>	 <p>Cloud Rendering</p>  <p>Cloud XR</p>  <p>Omniverse</p> <hr/> <p>Highest Perf Graphics Visual Computing</p> <p>Fastest RT Graphics Largest render models</p>	 <p>Virtual Desktop</p>  <p>Transcoding</p> <hr/> <p>Highest Density Virtual Desktop</p> <p>4K Resolution Max # of encode/decode streams</p>
<p>A100 250W & 300W 40G & 80G 2-slot FHFL NVLINK</p>	<p>A30 165W 24GB 2-slot FHFL NVLINK</p>	<p>T4 70W 16GB 1-slot Low Profile</p>	<p>A10 150W 24GB 1-slot FHFL</p>	<p>A40 300W 48GB 2-slot FHFL NVLINK</p>	<p>A16 250W 4 x 16GB 2-slot FHFL</p>

Compute

Graphics

NVIDIA A30

Versatile Compute Acceleration for Mainstream Enterprise Servers

Purpose built for Inference and Flexible Enterprise Compute

20X T4 AI perf (A30 TF32 FLOPS vs T4 FP32)

Multi-Instance GPU

Up to 4 concurrent instances per GPU (QoS)

Compute

3rd Gen Tensor cores, Fast FP64

High Bandwidth Memory

Ultra-low latency

Power Efficient

Excellent Perf/W

Sparsity Acceleration

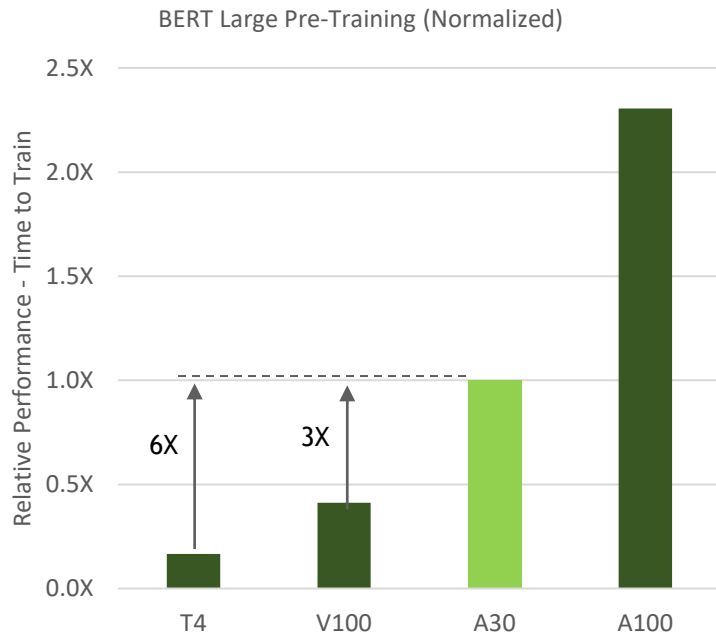
Further 2x speed up



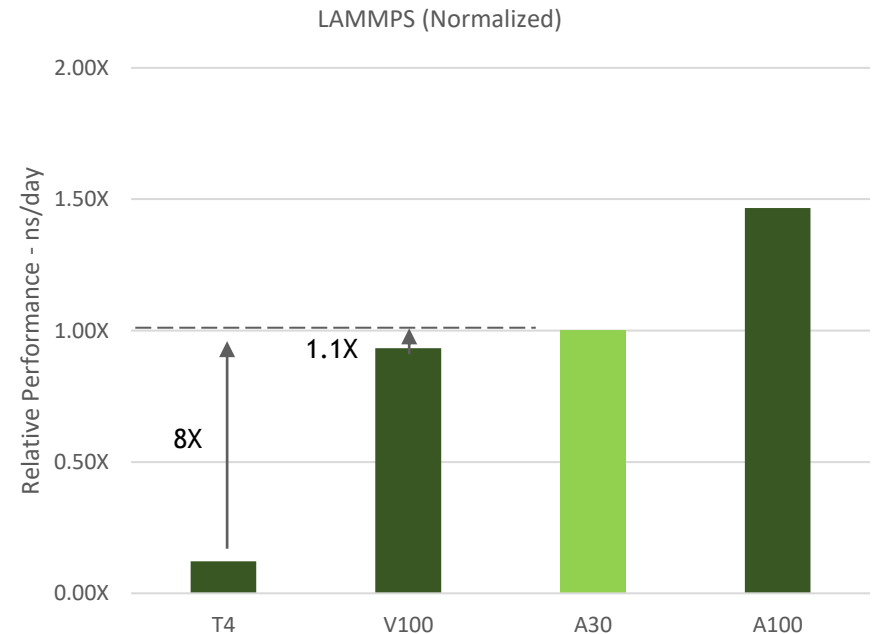
	A30
GPU Architecture	NVIDIA Ampere
Multi-Instance GPU	4 instances @ 6GB each 2 instances @ 12GB each
GPU Memory	24GB HBM2
Memory Bandwidth	933 GB/s
Interconnect	PCIe Gen 4 (x16) 1x NVLINK Bridge
Form Factor	2 Slot FHFL
Max Power	165W
Schedule	Production

A30 DELIVERED APPLICATION PERFORMANCE - TRAINING AND HPC

AI TRAINING—UP TO 3X HIGHER THROUGHPUT THAN V100 AND 6X HIGHER THAN T4



HPC— 1.1X HIGHER THROUGHPUT THAN V100 AND 8X HIGHER THAN T4

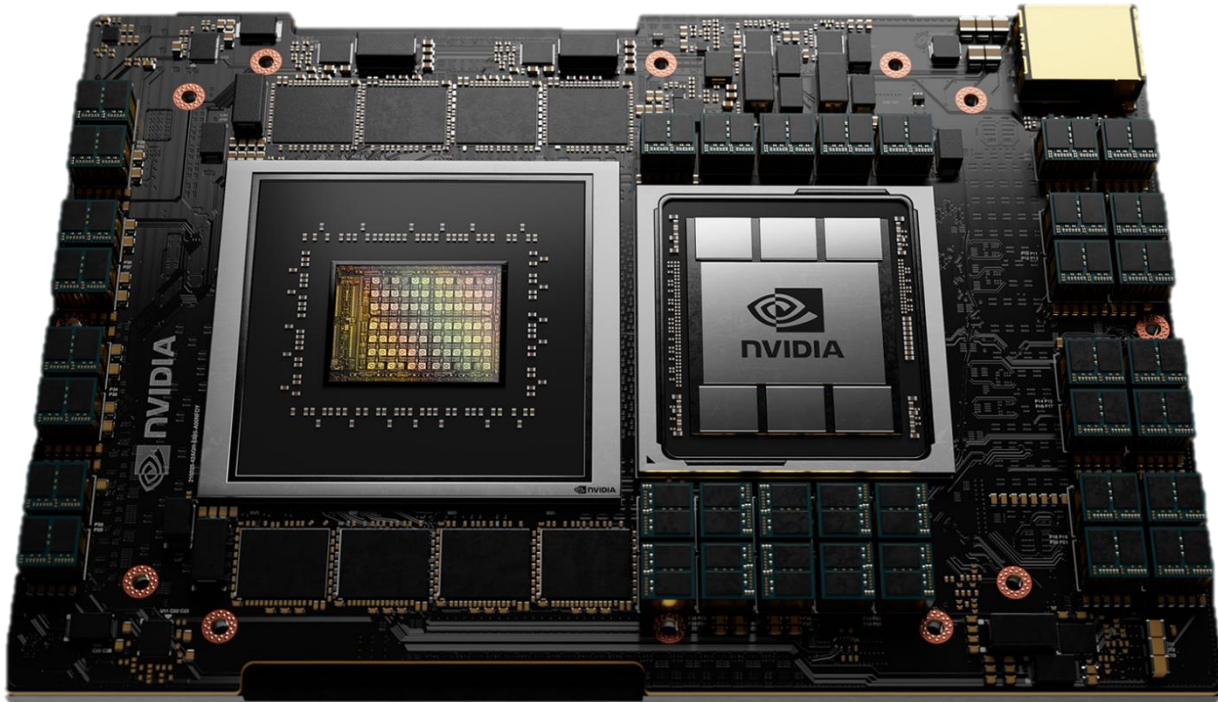


BERT-Large Pre-Training (9/10 epochs) Phase 1 and (1/10 epochs) Phase 2, Sequence Length for Phase 1 = 128 and Phase 2 = 512, dataset = real, NGC™ container = 21.03, 8x GPU: T4 (FP32, BS=8, 2) | V100 PCIE 16GB (FP32, BS=8, 2) | A30 (TF32, BS=8, 2) | A100 PCIE 40GB (TF32, BS=54, 8) | batch sizes indicated are for Phase 1 and Phase 2 respectively

Dataset: ReaxFF/C, FP64 | 4x GPU: T4, V100 PCIE 16GB, A30

ANNOUNCING NVIDIA GRACE

Breakthrough CPU Designed for Giant-Scale AI and HPC Applications



FASTEST INTERCONNECTS

>900 GB/s Cache Coherent NVLink CPU To GPU (14x)
>600GB/s CPU To CPU (2x)

HIGHEST MEMORY BANDWIDTH

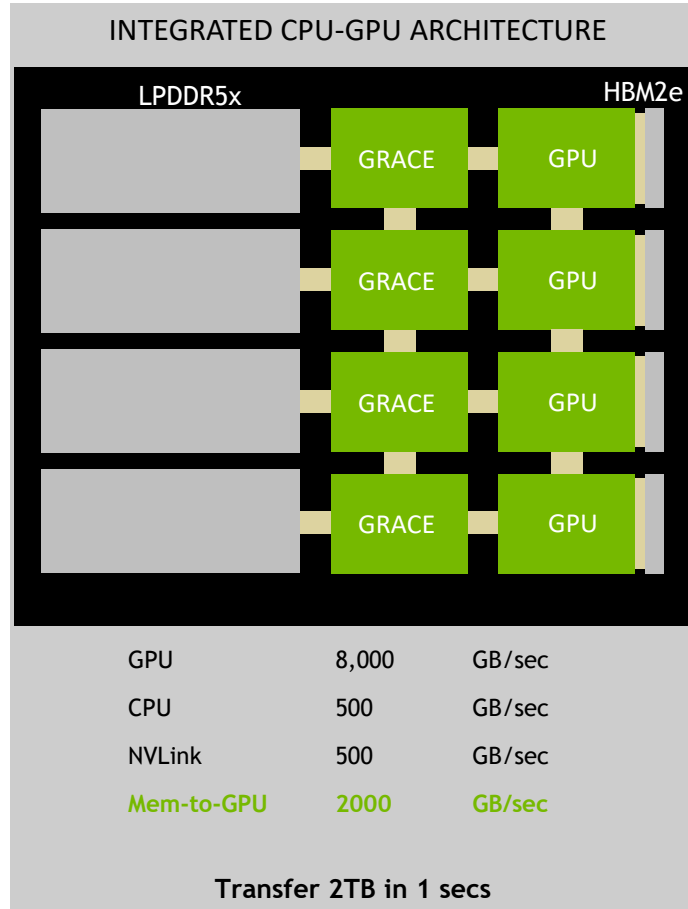
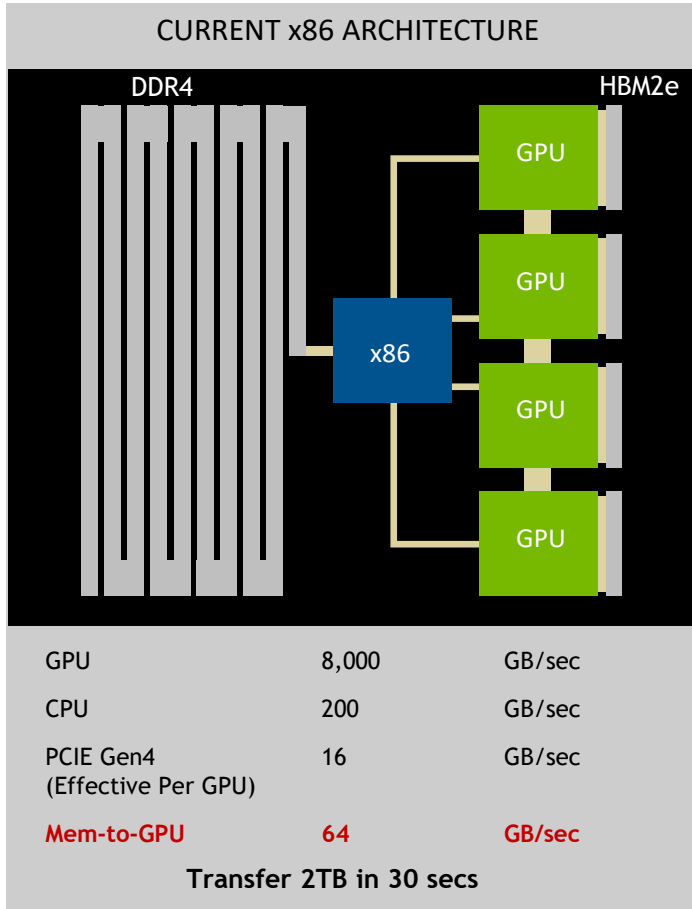
>500GB/s LPDDR5x w/ ECC
>2x Higher B/W
10x Higher Energy Efficiency

NEXT GENERATION ARM NEOVERSE CORES

>300 SPECrate2017_int_base est.
Availability 2023

TURBOCHARGED TERABYTE SCALE ACCELERATED COMPUTING

Evolving Architecture For New Workloads



3 DAYS FROM 1 MONTH
Fine-Tune Training of 1T Model

**REAL-TIME INFERENCE
ON 0.5T MODEL**
Interactive Single Node NLP Inference

Bandwidth claims rounded to nearest hundred for illustration.
Performance results based on projections on these configurations Grace : 8xGrace and 8xA100 with 4th Gen NVIDIA NVLink Connection between CPU and GPU and x86: DGX A100.
Training: 1 Month of training is Fine-Tuning a 1T parameter model on a large custom data set on 64xGrace+64xA100 compared to 8xDGX A100 (16xX86+64xA100)
Inference: 530B Parameter model on 8xGrace+8xA100 compared to DGXA100.

ANNOUNCING THE WORLD'S FASTEST SUPERCOMPUTER FOR AI

20 Exaflops of AI

Accelerated w/ **NVIDIA Grace CPU and NVIDIA GPU**

HPC and AI For Scientific and Commercial Apps

Advance Weather, Climate, and Material Science





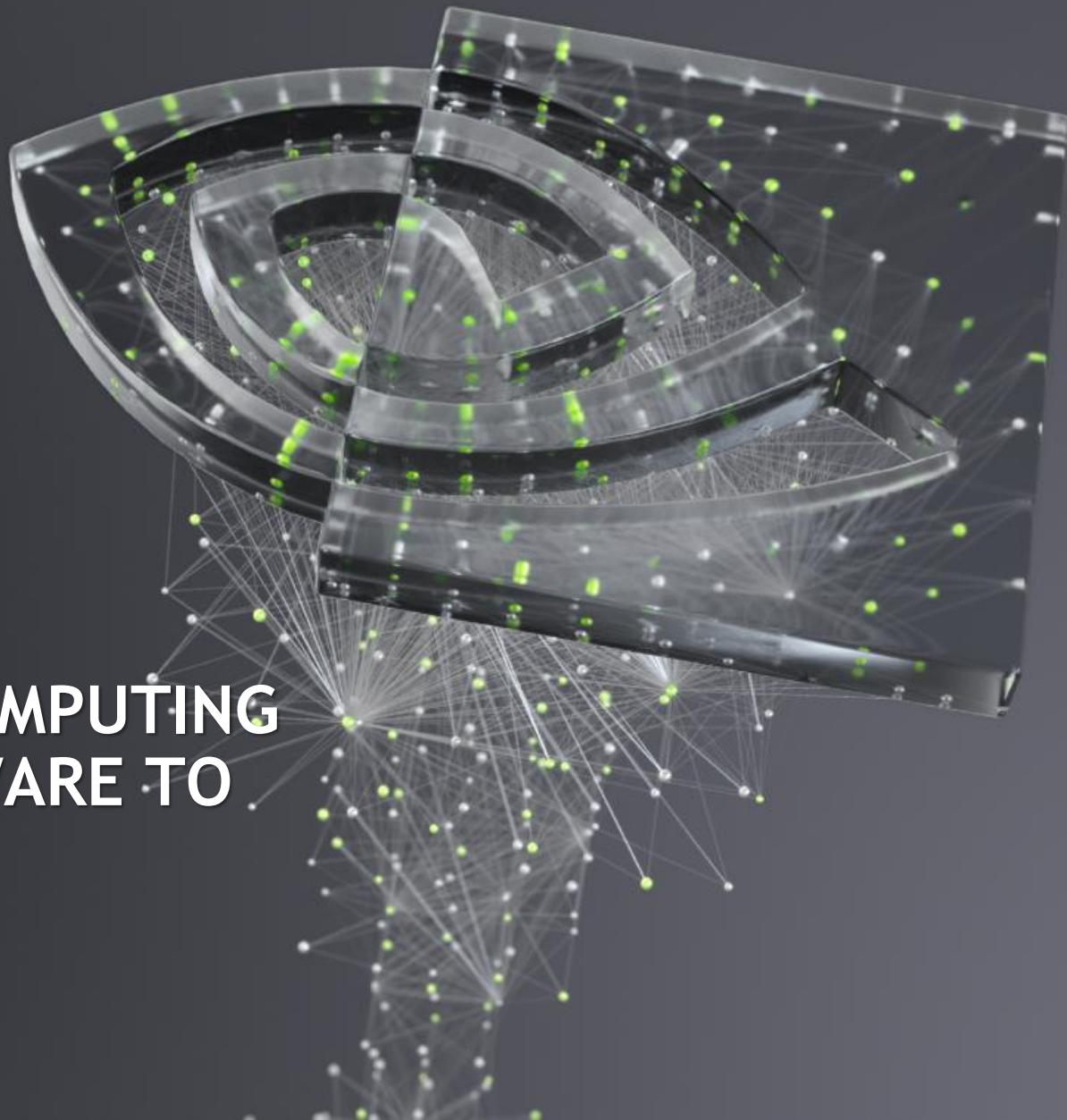
TECH TAKE

NVIDIA INSIGHTS: GPU COMPUTING LANDSCAPE FROM SOFTWARE TO HARDWARE.

Gabriel Noaje, PhD

Senior Solutions Architect, NVIDIA APAC South

June 29, 2021



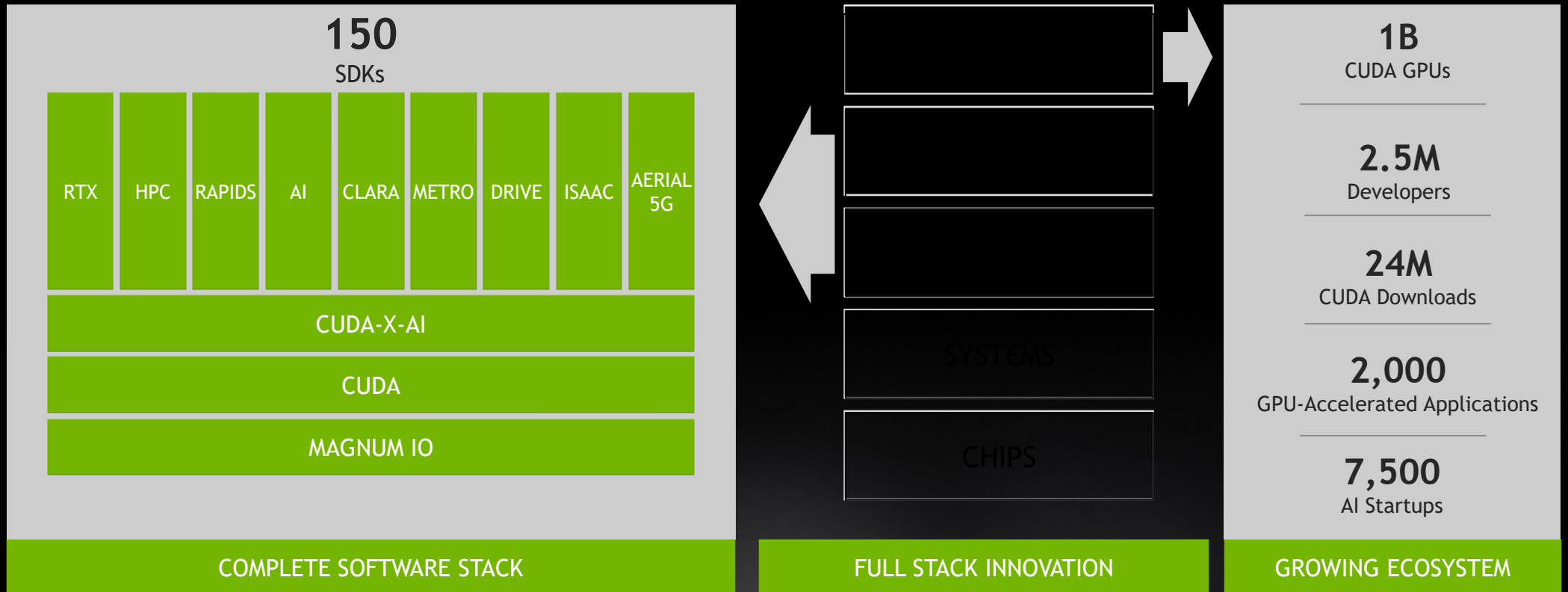
NVIDIA IS A FULL STACK COMPUTING PLATFORM

Amazing Innovation and Expansion of NVIDIA Ecosystem



NVIDIA IS A FULL STACK COMPUTING PLATFORM

Amazing Innovation and Expansion of NVIDIA Ecosystem



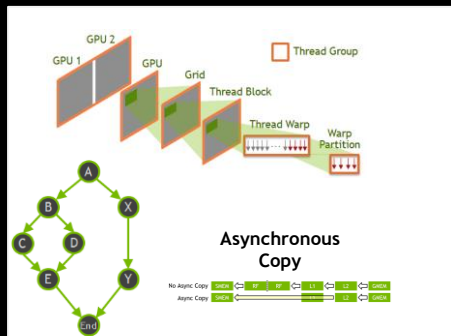
CUDA 11.4

Major Feature Areas



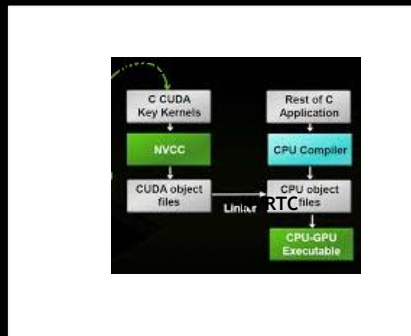
New Platform Capabilities

-
-



Programming Model

- *cudaGraphs*
- *cudaMallocAsync*
- *Driver symbols*
- *Language support:*
 - *CUDA Python*



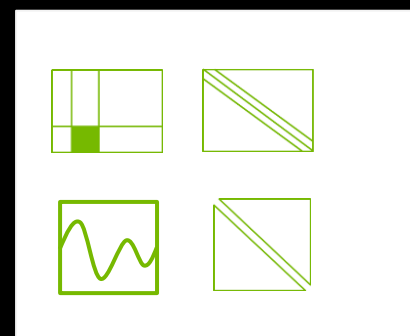
Compiler

- *cudaMallocAsync* builtin-function
- *alloca()* builtin-function



Developer Tools

- *cudaGraphs*
- *GPU metrics sampling*



Math Libraries

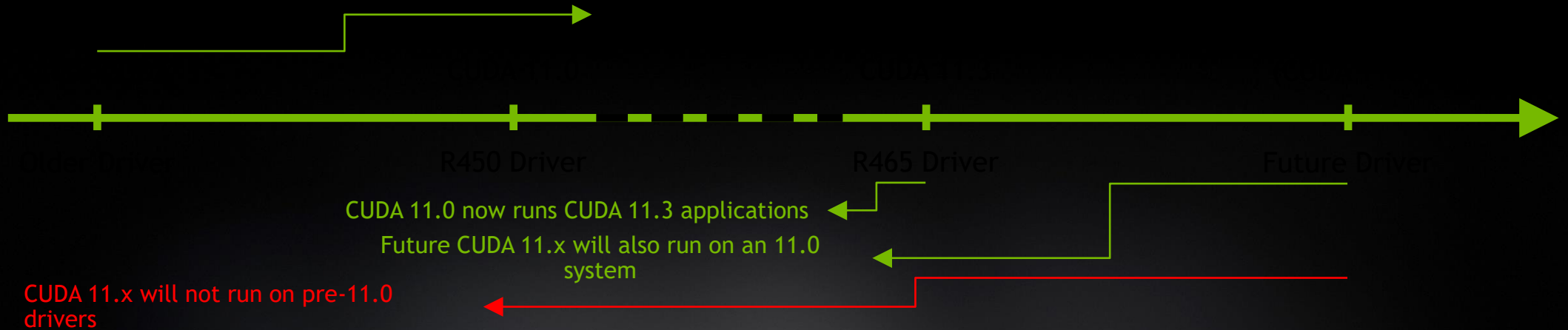
-

MINOR VERSION COMPATIBILITY

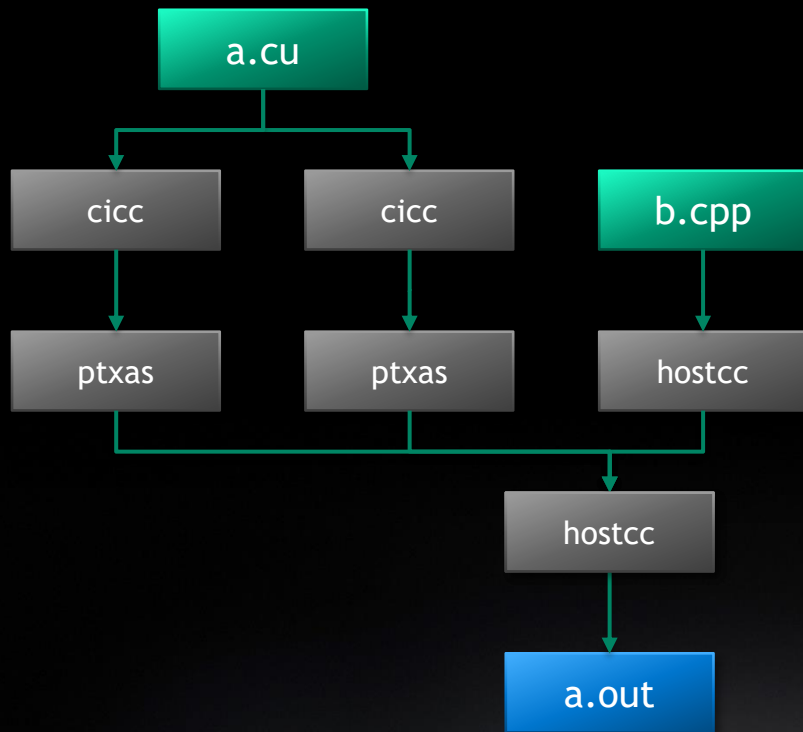
always

newer

same major version



NVCC MULTI-TARGET PARALLEL COMPILATION

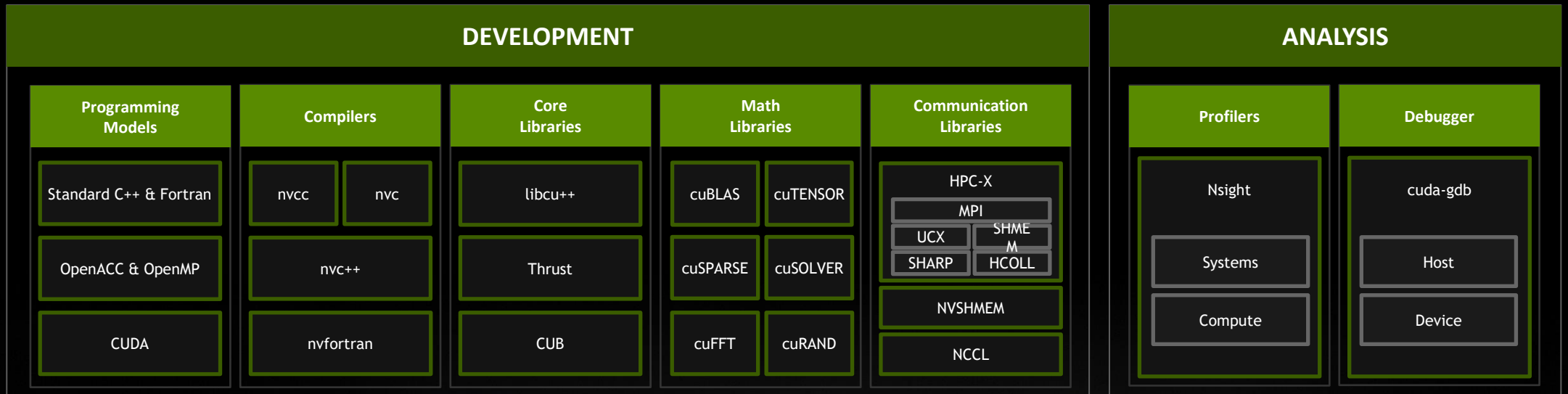


Parallel compilation for multiple architectures

```
# nvcc -t x.cu y.cu z.cu ← Multi-target, single arch  
  
# nvcc -t a.cu b.cpp -gencode=arch=compute_70,code=sm_70  
\  
-gencode=arch=compute_80,code=sm_80
```

NVIDIA HPC SDK

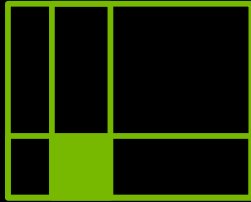
Now including NVIDIA HPC-X Communications Stack



Develop for the NVIDIA Platform: GPU, CPU and Interconnect
Libraries | Accelerated C++ and Fortran | Directives | CUDA
7-8 Releases Per Year | Freely Available

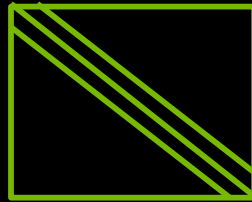
GPU ACCELERATED MATH LIBRARIES IN CUDA 11

OVERVIEW OF NEW NVIDIA A100 FEATURES



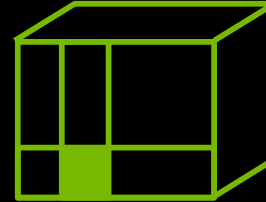
cuBLAS

BF16, TF32 and
FP64 Tensor
Cores



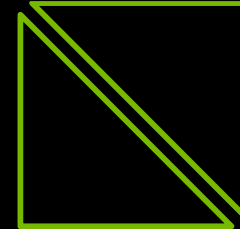
cuSPARSE

BF16, Increased
memory BW, Shared
Memory and L2



cuTENSOR

BF16, TF32 and
FP64 Tensor
Cores



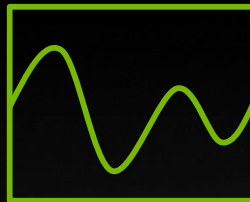
cuSOLVER

BF16, TF32 and
FP64 Tensor
Cores



CUTLASS

BF16, TF32 and
FP64 Tensor
Cores



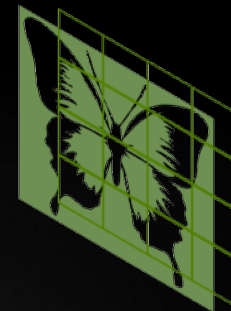
cuFFT

Increased memory
BW, Shared Memory
and L2



CUDA Math API

BF16 Support



nvJPEG

Hardware
Decoder

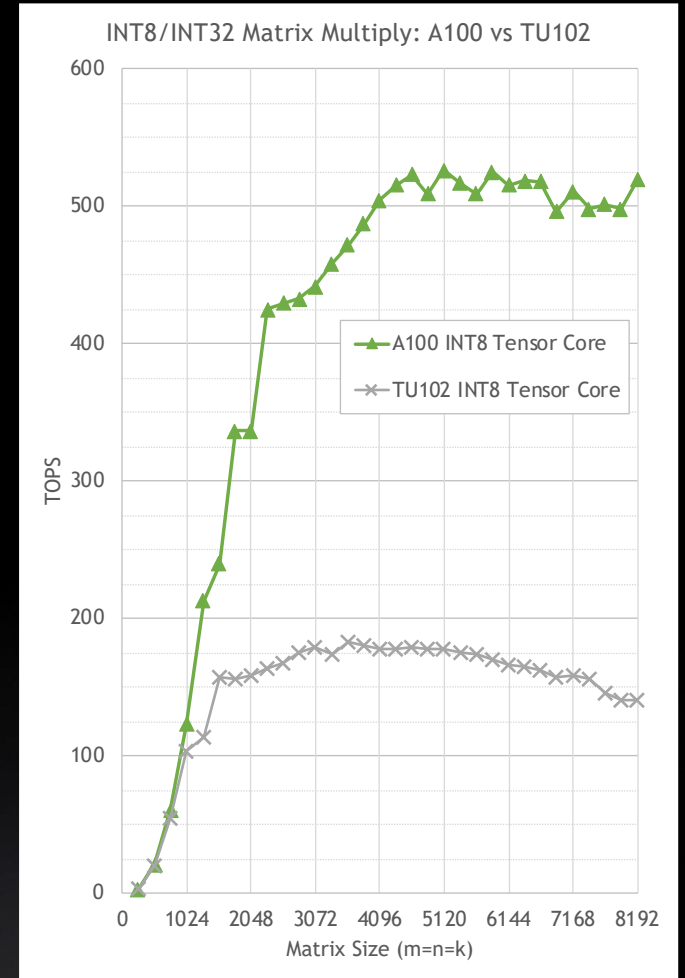
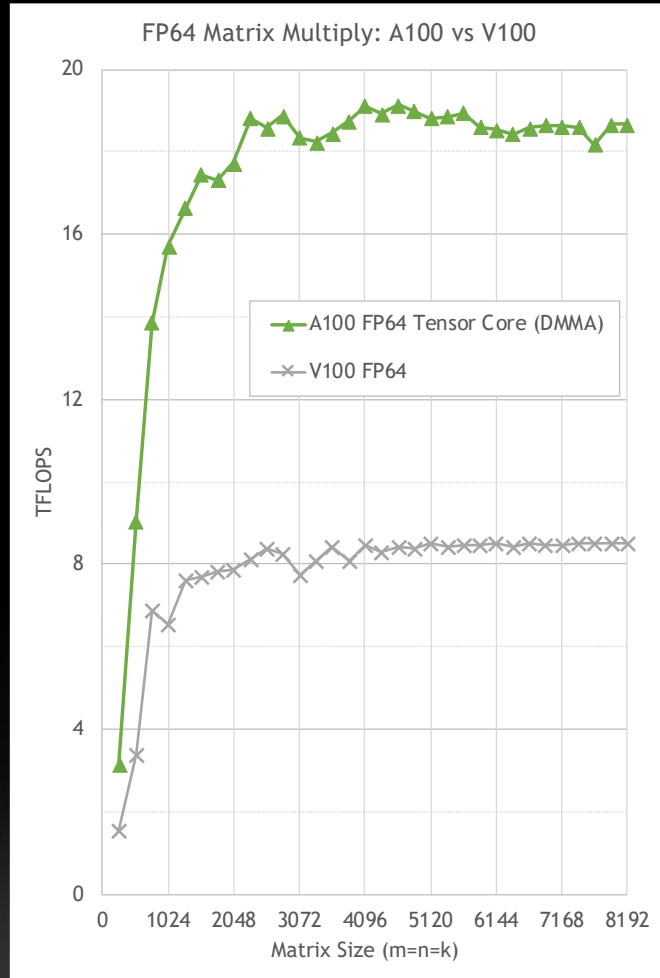
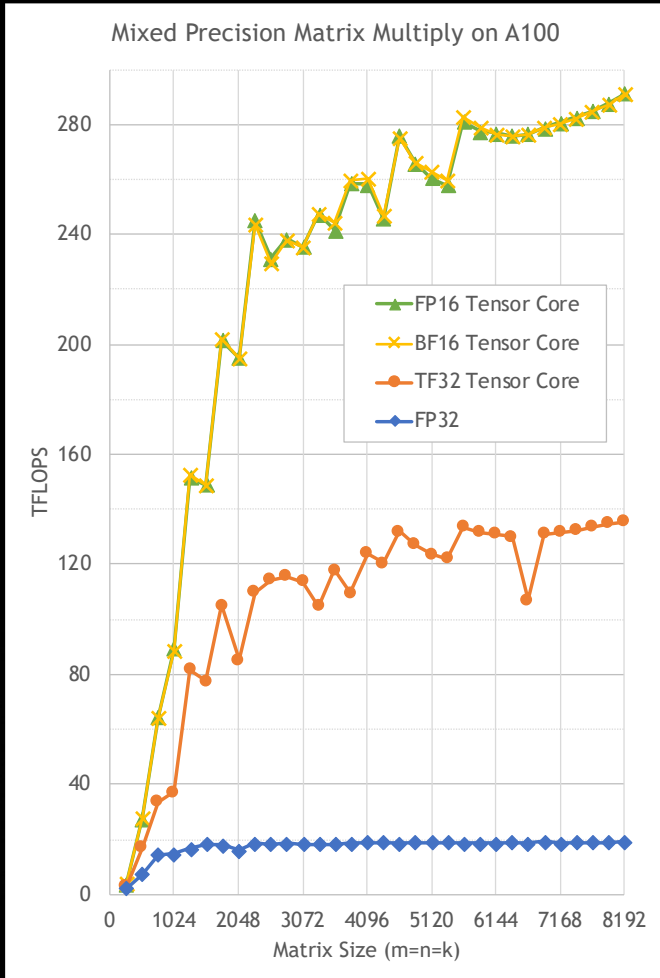
TENSOR CORE SUPPORT IN MATH LIBRARIES

High-level overview of supported functionality by each library

	INT4		INT8		FP16		BF16		TF32		FP64
	Dense	Sparse	Dense	Sparse	Dense	Sparse	Dense	Sparse	Dense	Sparse	Dense
cuBLAS			✓		✓		✓		✓		✓
cuBLAS Sparse					✓		✓		✓		✓
cuBLAS Sparse					✓		✓		✓		✓
cuBLAS Sparse			✓		✓		✓		✓		✓
CUTLASS SpMM				✓		✓		✓		✓	
CUTLASS Dense GEMM and SpMM	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
CUTLASS Convolutions	✓		✓		✓		✓		✓		

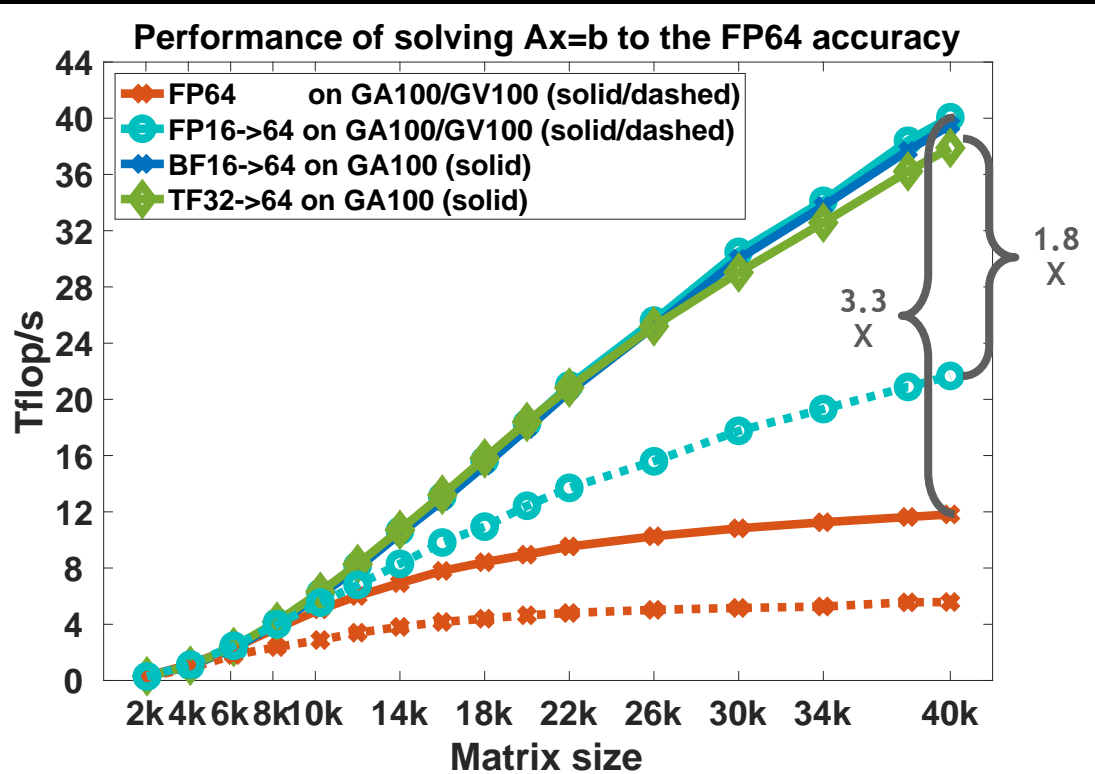
cuBLAS

3rd GENERATION TENSOR CORES ADD SUPPORT FOR FP64 & NEW TYPE BF16 & COMPUTE TYPE TF32

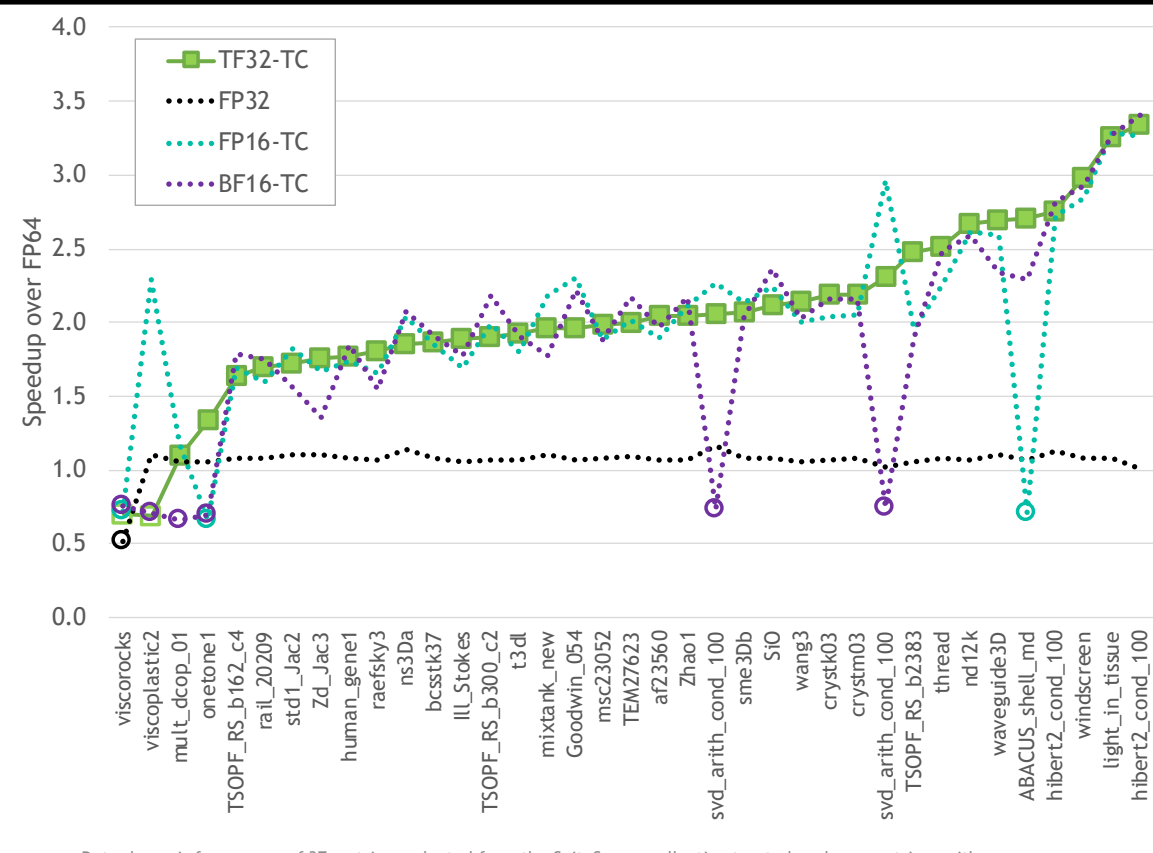


cuSOLVER

TENSOR CORE ACCELERATED ITERATIVE REFINEMENT SOLVER (TCAIRS)



Matrices used in this performance chart are Hilbert matrices. TCAIRS solver is ~3.3X faster than full FP64 solver using DMMA tensor cores and A100 is ~1.8X faster than V100.

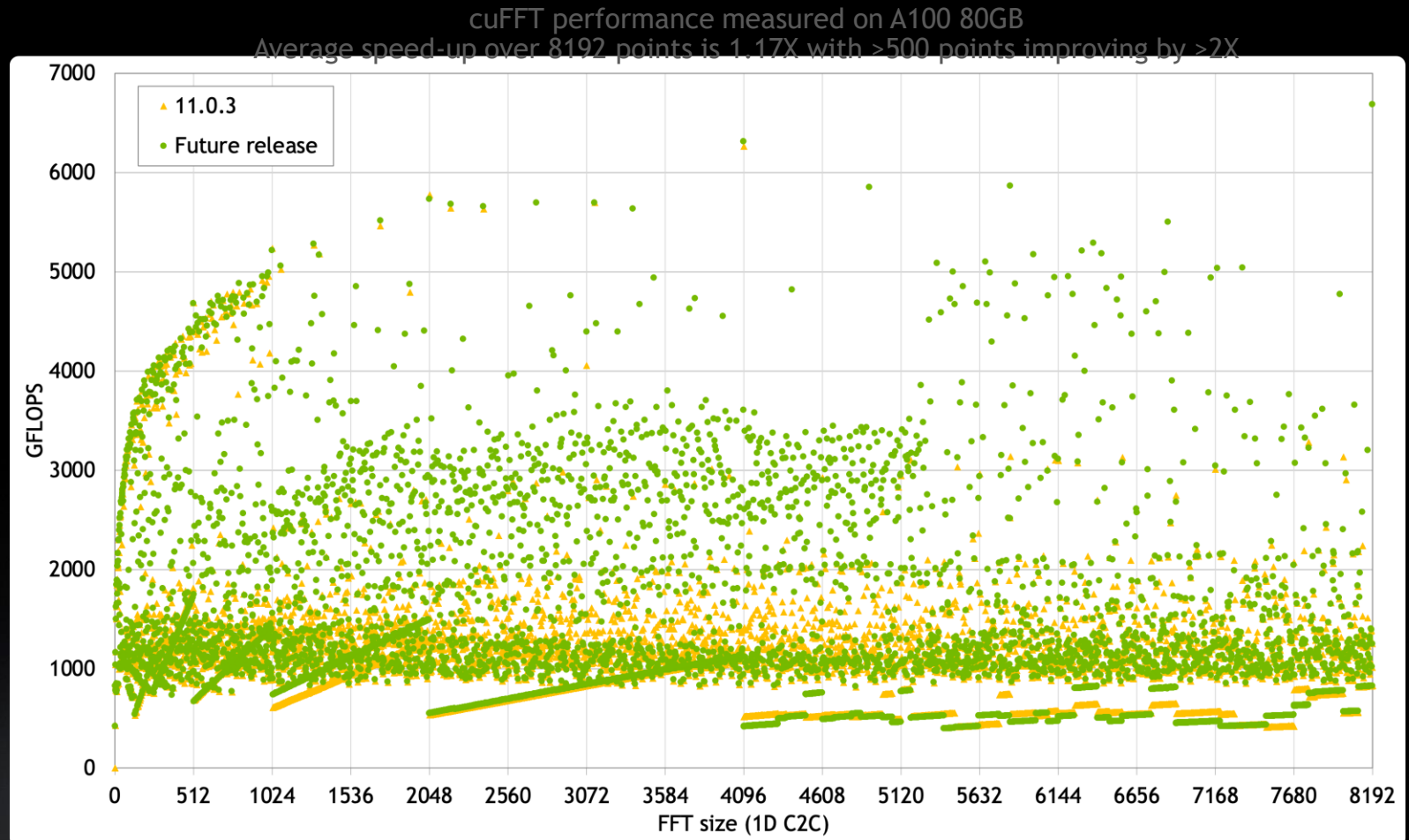


Data shown is for a group of 37 matrices selected from the SuiteSparse collection treated as dense matrices with sizes ranging from 20k to 38k. Performance data is average of all cases and all runs were done on an NVIDIA A100 with cuSOLVER from CUDA Toolkit 11.0. FP16 cases were scaled to ensure matrix entries are within the type range.

IMPROVED PERFORMANCE cuFFT

Updated performance and features since CUDA 11.0

- ▶ Optimized FFTs up to size 512 per dimension
- ▶ Improved C2C, R2C and C2R transforms
- ▶ bfloat16 support since 11.1.1
- ▶ Coming soon:
 - ▶ Optimized FFT sizes up to size 32768 per dimension
 - ▶ Improved decompositions for large FFT sizes (>16384) and 1D/2D/3D



DEVICE EXTENSION LIBRARIES

Enabling kernel fusion of high-performance numerical method implementations

▶ cuFFTDx: In Math Library EA Program

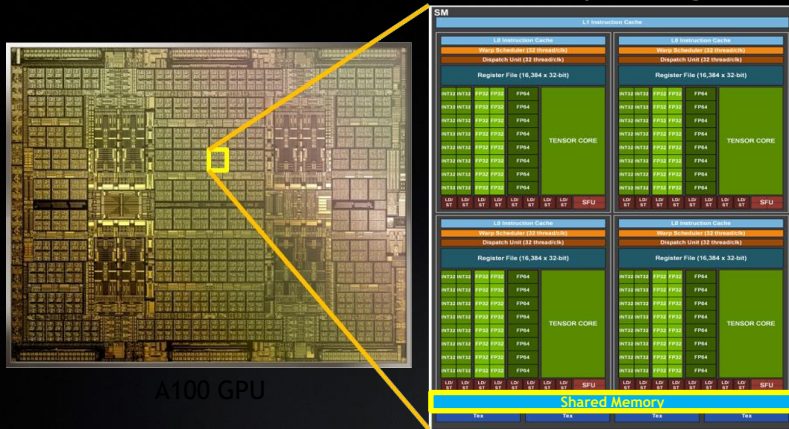
- ▶ Retain and reuse of on-chip data
- ▶ Inline FFTs in user kernel with sizes up to 32k (on A100)
- ▶ Combine FFT operations
- ▶ Expecting cuFFTDx GA release later this year

```
using namespace cufftdx;

using FFT = decltype(Block() + Size<128>() + Type<fft_type::c2c>() +
    Direction<fft_direction::forward>() + Precision<float>() +
    ElementsPerThread<8>() + FFTsPerBlock<2>() + SM<700>());
```

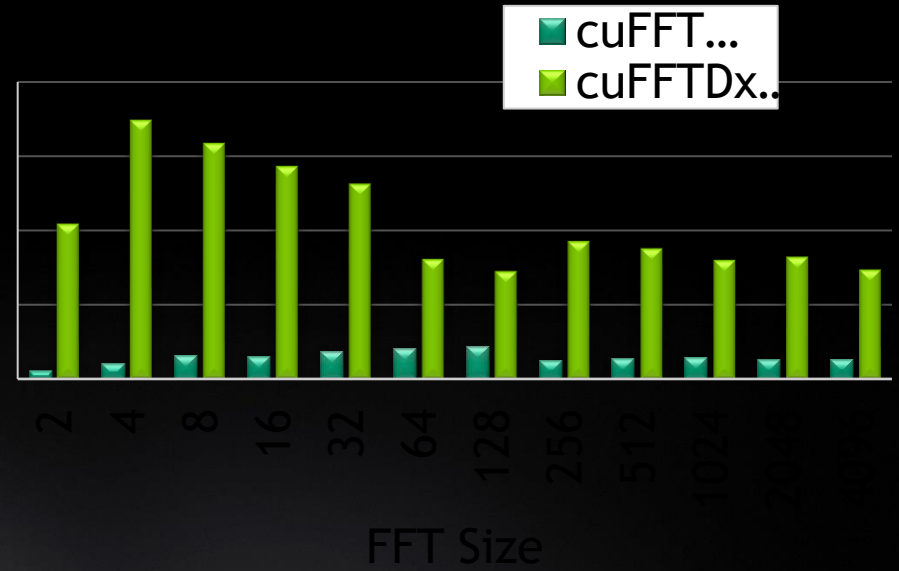
▶ Linear algebra device extensions

- ▶ cuSOLVERDx with non-pivoting LU also in EA
- ▶ cuBLASDx with GEMM and TRSM are upcoming



A100 GPU

A100 Streaming Multiprocessor



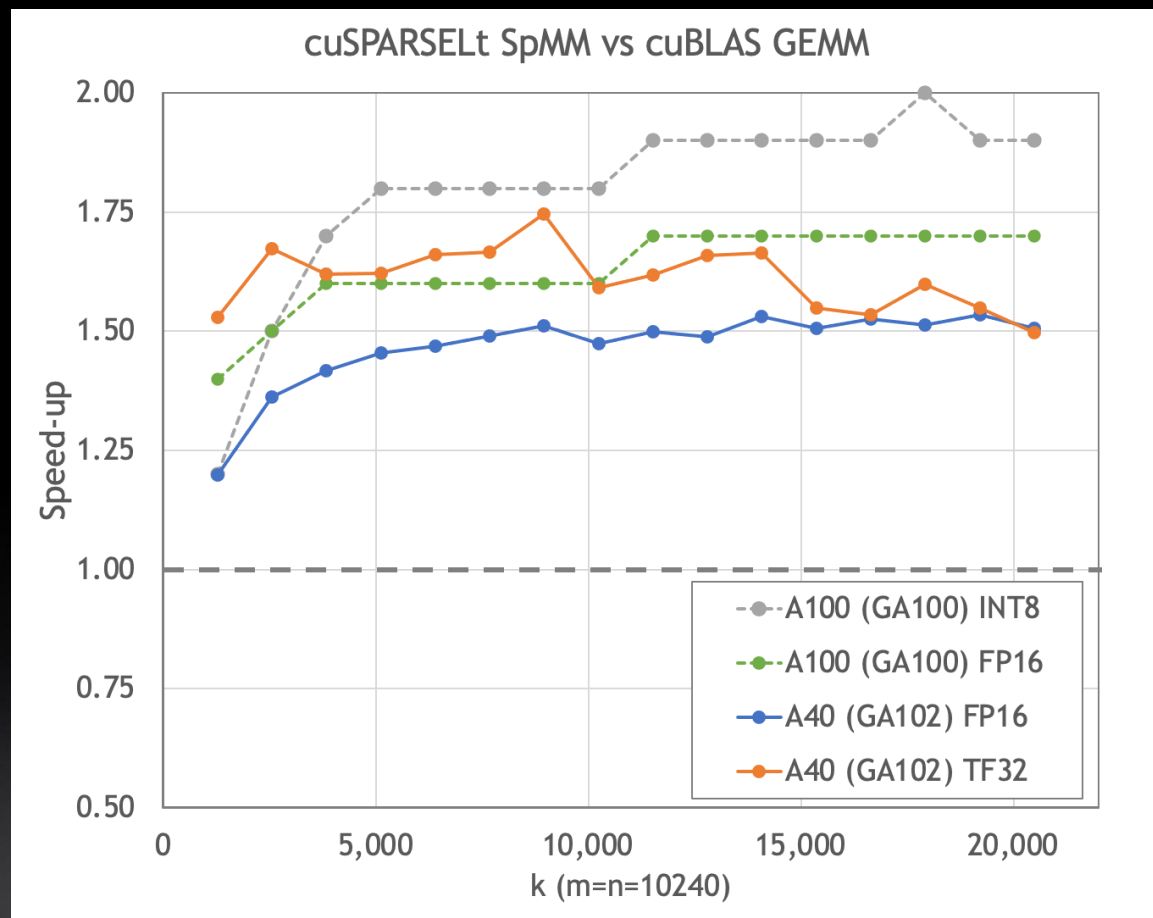
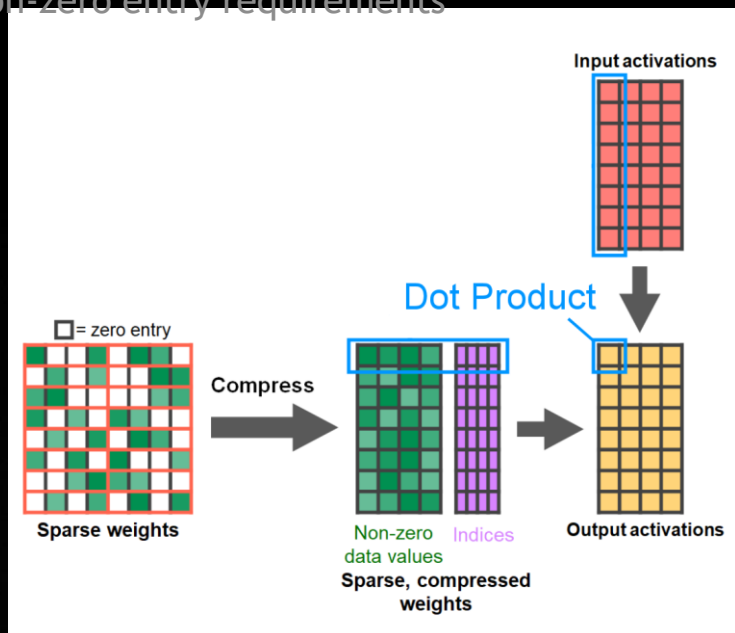
<https://developer.nvidia.com/CUDAMathLibraryE>

A

EXPLOITING STRUCTURED SPARSITY WITH cuSPARSELt ON NVIDIA AMPERE ARCHITECTURE GPUS

Most recent update is available at <https://developer.nvidia.com/cusparselt/downloads>

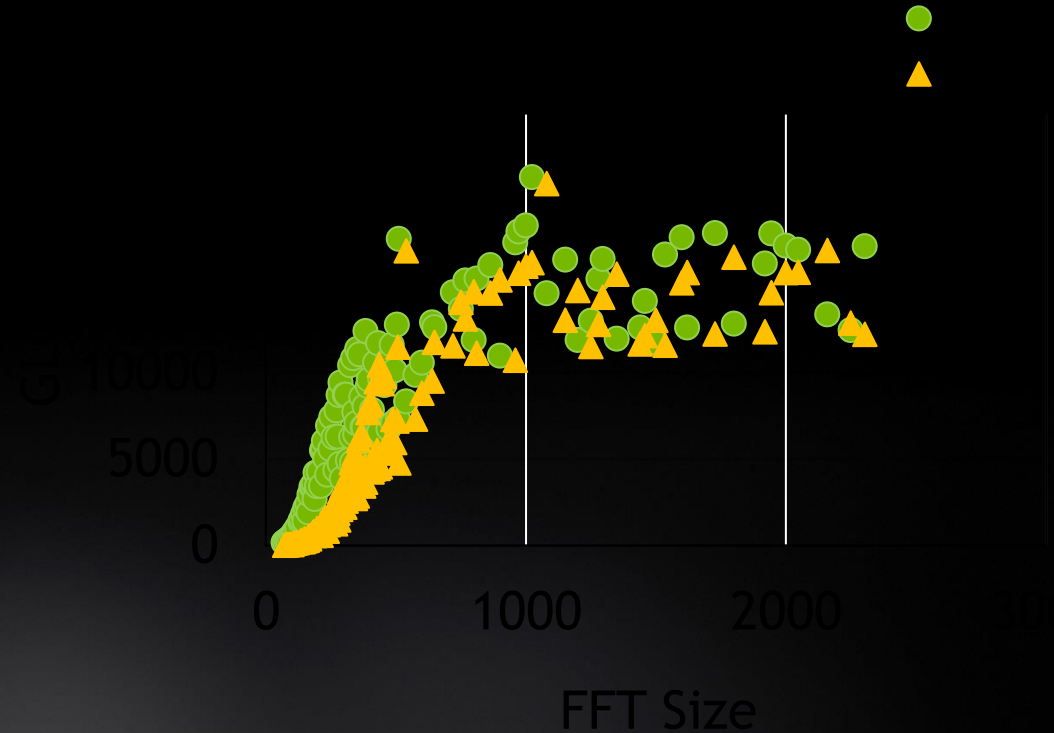
- ▶ Provides easy access to 3rd generation Tensor Core Sparse MMAs on NVIDIA Ampere architecture GPUs to perform Sparse x Dense Matrix Multiplies
- ▶ cuSPARSELt provides essential utilities to easily create 50% sparse matrices that conform to the 2:4 non-zero entry requirements



MULTI-GPU cuFFT

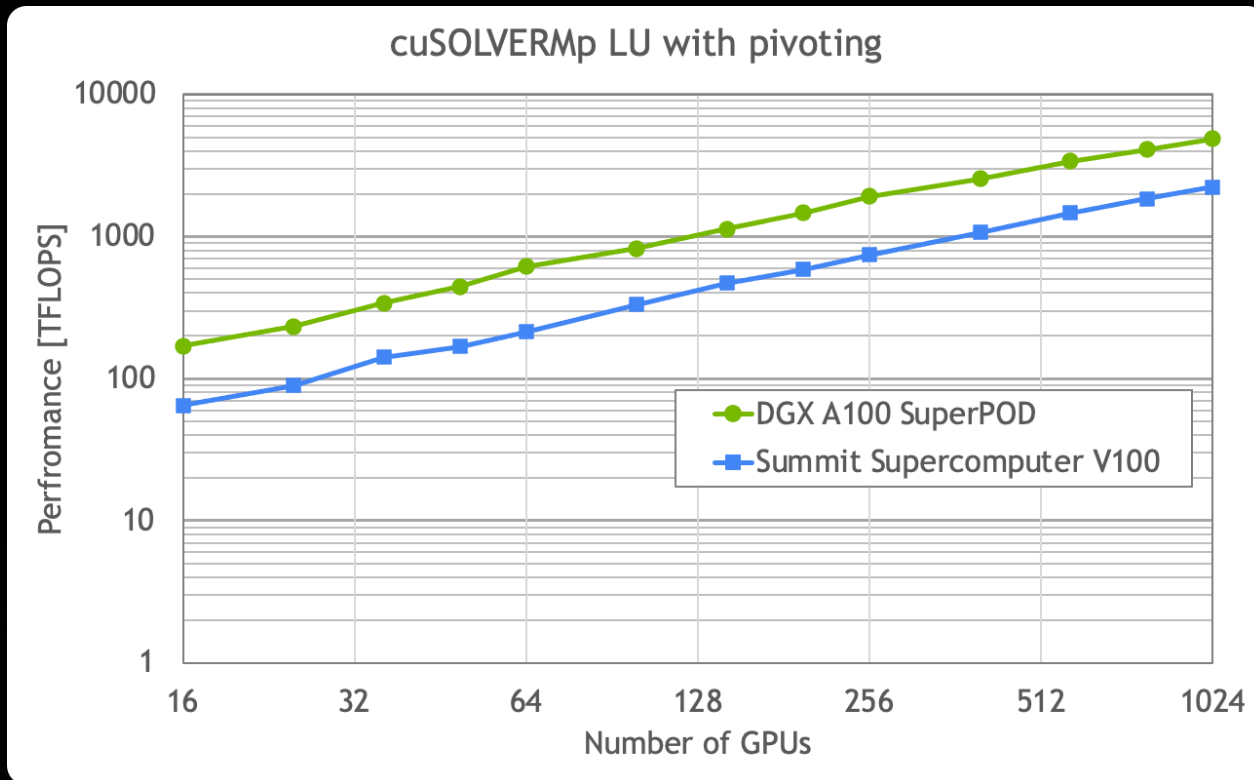
Single-node multi-GPU improvements available with CUDA 11.2

- ▶ Improved strong scaling performance for small 2D/3D FFTs with size <512
- ▶ Up to 22 TFLOPS SP & 11 TFFLOPS DP of performance on a DGX A100 node for 2D/3D FFT transforms
- ▶ Maximum size per node increased to 3072 cubed 3D FFT (SP)
- ▶ Stream ordering/graph capturing added to multi-GPU plans
- ▶ Upcoming later this year: Multi-process API



cuSOLVER DISTRIBUTED MULTI-GPU LU SOLVER

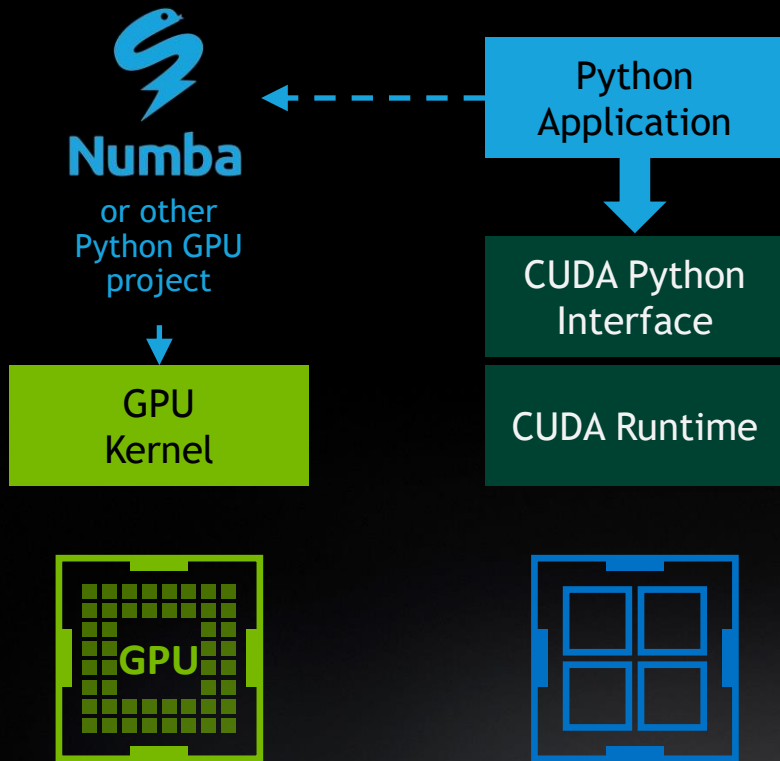
Scalability on Top 5 supercomputers



Data shown is for double complex LU factorization for 27882 rows/cols per device

- ▶ Supports standard 2D block cycling distribution of matrix and rhs between processes
- ▶ > 2.2X speed-up going from V100 to A100
- ▶ First release is planned to be part of HPC SDK in 1st half of 2021
- ▶ Early Access available by end of April
<https://developer.nvidia.com/cudamathlibrarye>
[a](#)
- ▶ Further performance improves to come after initial release

CUDA PYTHON



Distribution

- Source available on GitHub
- PIP & Conda packages
- Redistributable license

Bindings

Full coverage of and access to the CUDA host APIs from Python

Platforms

- Linux: `x86_64`, `sbsa`, `ppc64le`
- Windows: `x86_64`

INTRODUCING LEGATE

Accelerated and Distributed

A framework for programming large numbers of GPUs as if they were a single processor

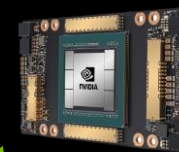
Pass data between Legate libraries without worrying about distribution or synchronization requirements

Legate NumPy and Pandas aim to transparently scale existing Numpy and Pandas workloads

Legate Numpy and Legate Pandas available now and opensource!

```
import legate.numpy as np

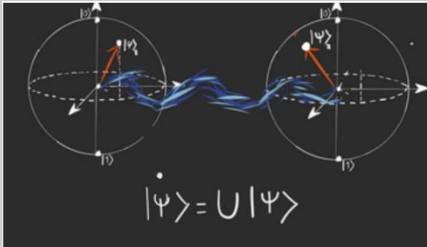
def cg_solve(A, b, tol=1e-10):
    x = np.zeros(A.shape[1])
    r = b - A.dot(x)
    p = r
    rsold = r.dot(r)
    for i in xrange(b.shape[0]):
        Ap = A.dot(p)
        alpha = rsold / (p.dot(Ap))
        x = x + alpha * p
        r = r - alpha * Ap
        rsnew = r.dot(r)
        if np.sqrt(rsnew) < tol:
            break
        beta = rsnew / rsold
        p = r + beta * p
        rsold = rsnew
    return x
```



ANNOUNCING NVIDIA CUQUANTUM

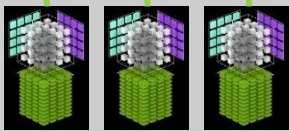
Researching the Computer of Tomorrow on the Most Powerful Computer Today

SDK FOR GPU-ACCELERATED QUANTUM SIMULATIONS

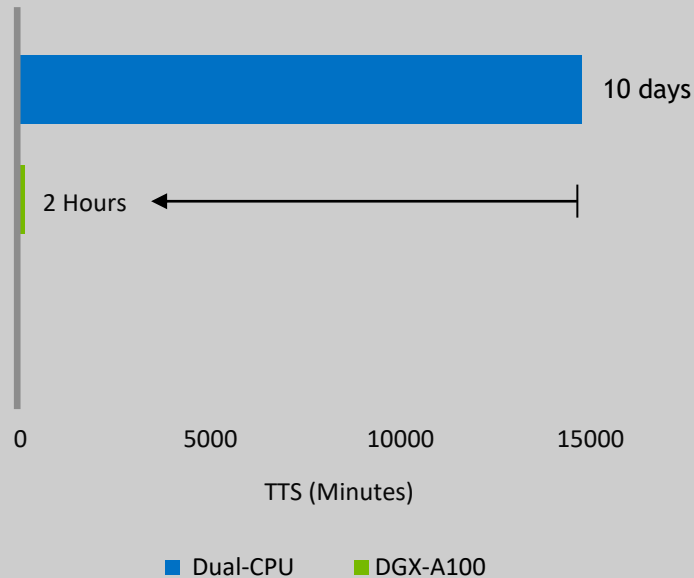


Quantum Circuit Simulation Frameworks

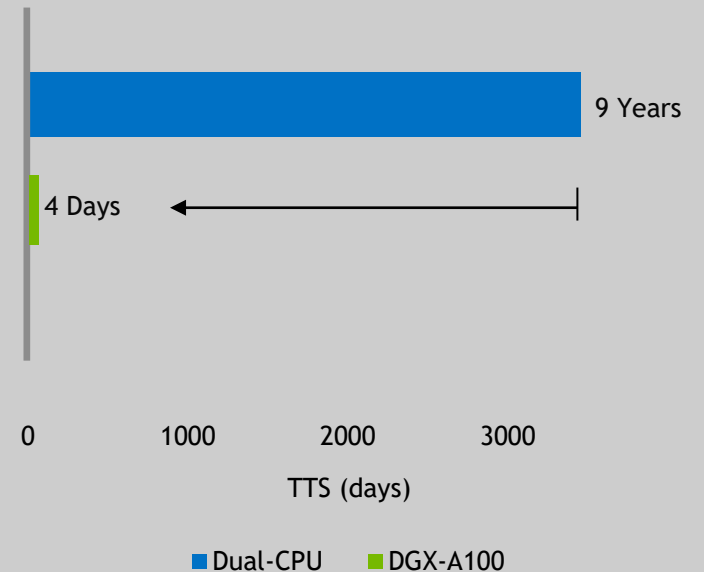
cuQuantum



STATE VECTOR SIMULATION Scales to 10s of Qubits



TENSOR NETWORK SIMULATION Scales to 1000's of Qubits



Data Analytics & Training

HugeCTR mxnet ONNX

PaddlePaddle PYTORCH

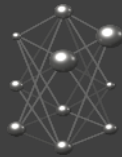
RAPIDS APACHE spark TensorFlow

CUDA-X-AI

CUDA

MAGNUM IO

Inference



Conversational AI
Computer Vision
Recommendations
Reinforcement Learning

Pre-Trained Models
SOTA Models



Triton
Inference Serving



Transfer Learning Toolkit
Zero Coding | Speech, Vision & NLU

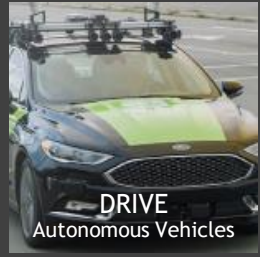


TensorRT 7.2

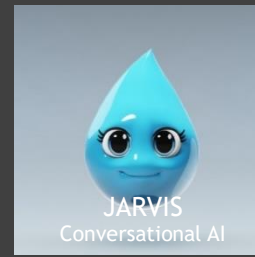
Application Frameworks



CLARA
Healthcare



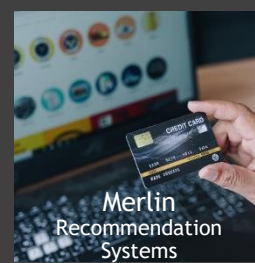
DRIVE
Autonomous Vehicles



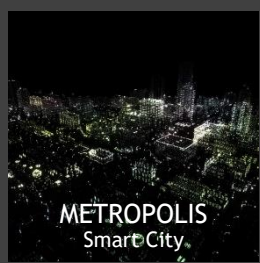
JARVIS
Conversational AI



ISAAC
Robotics

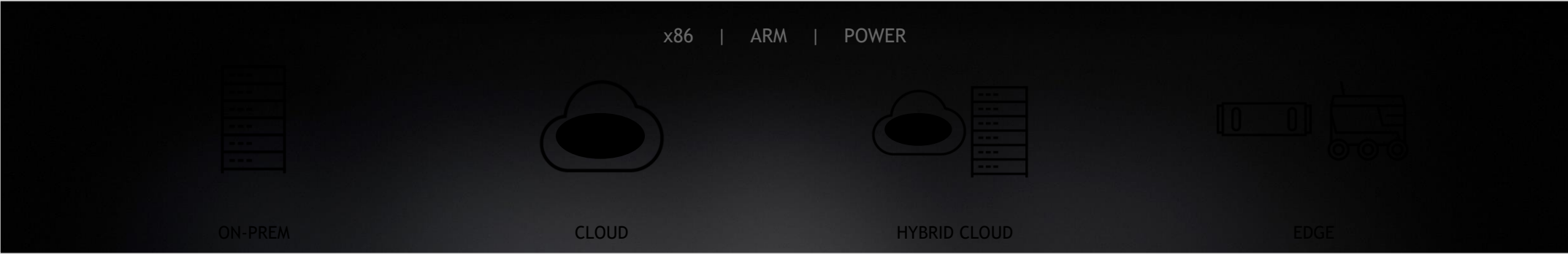
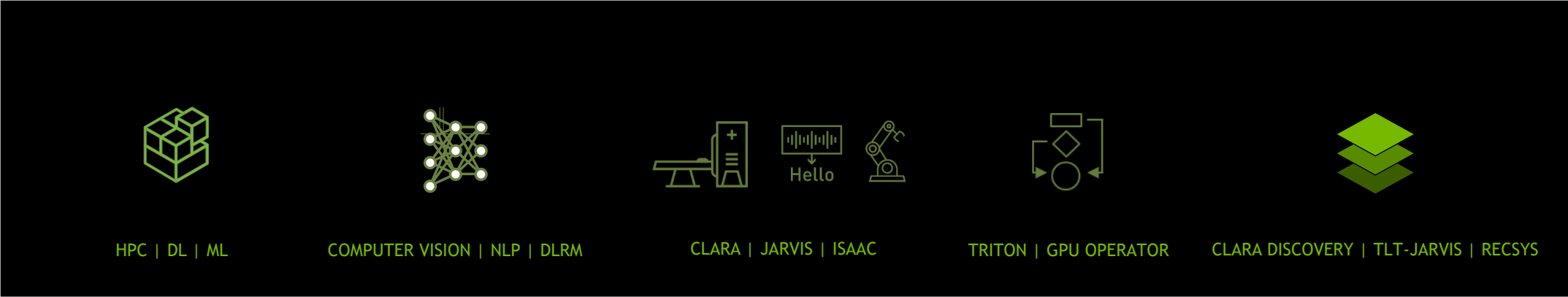


Merlin
Recommendation
Systems

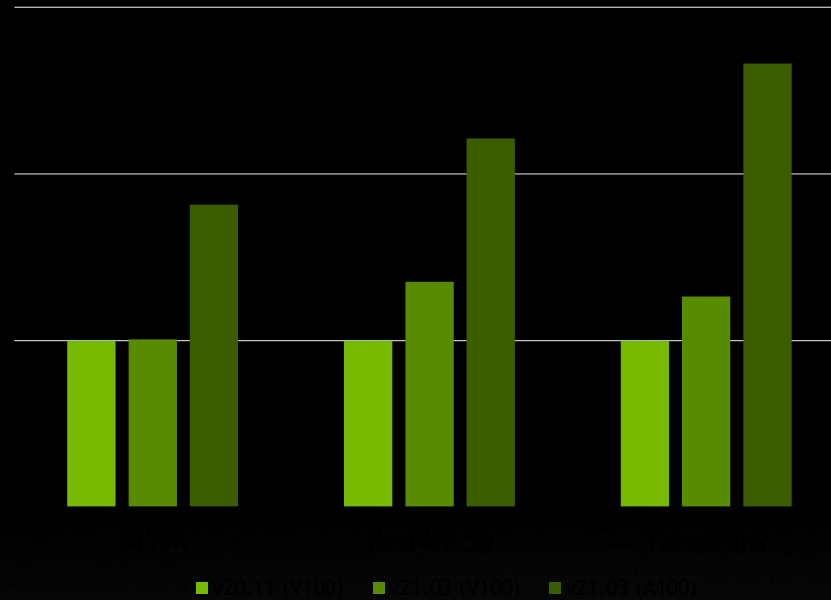


METROPOLIS
Smart City

Build AI Faster, Deploy Anywhere



DO WHAT YOU DO BEST, FASTER



ENTERPRISE READY SOFTWARE

- ✓ Tested for CVEs, malware, crypto
- ✓ Tested for reliability
- ✓ Backed by Enterprise support

PERFORMANCE OPTIMIZED

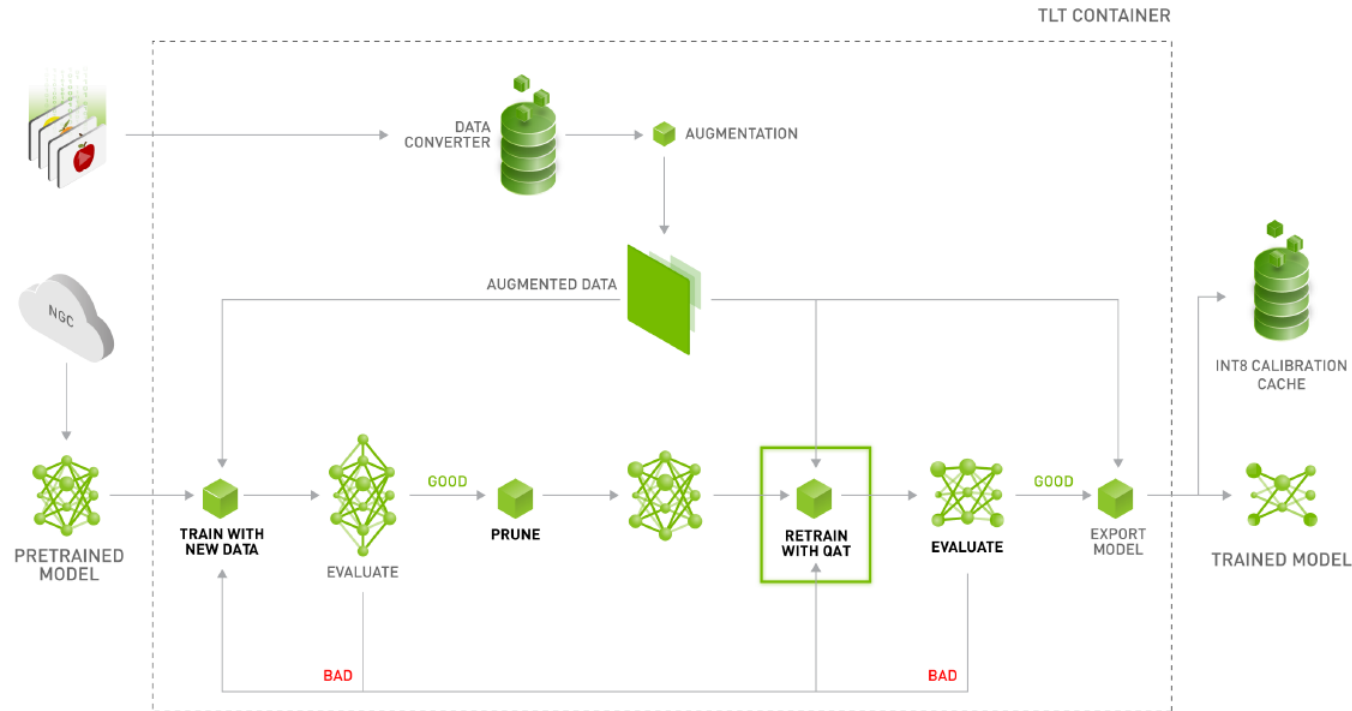
- ✓ Scalable
- ✓ Updated Monthly
- ✓ Better performance on the same system

DEPLOY ANYWHERE

- ✓ Docker | cri-o | containerd | OpenShift
- ✓ Bare metal, VMs, Kubernetes
- ✓ Multi-cloud, on-prem, hybrid

TLT WORKFLOW

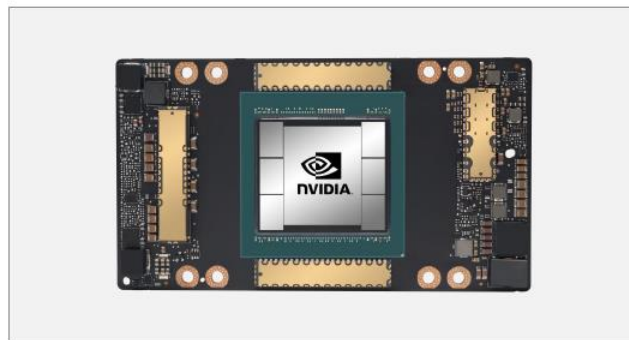
Automatic Mixed Precision | Quantization Aware Training | Pruning



ANNOUNCING TLT 3.0 DEVELOPER PREVIEW



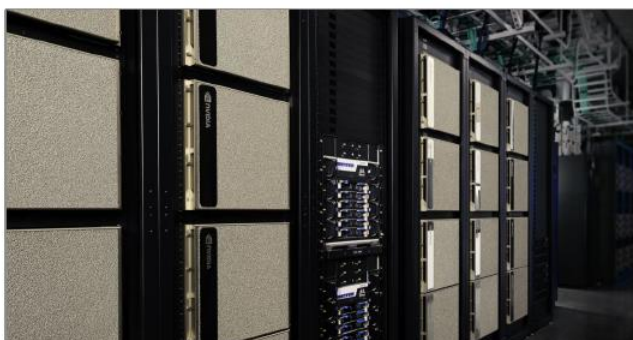
TRAIN CONVERSATIONAL AI MODELS
Simplify training of complex ASR & NLP tasks



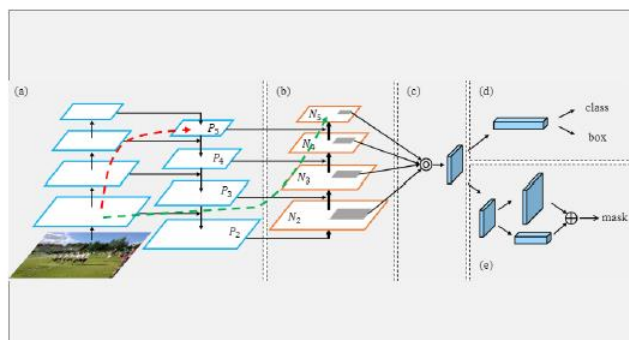
ACCELERATE TRAINING WITH AMPERE
10x speedup in training time



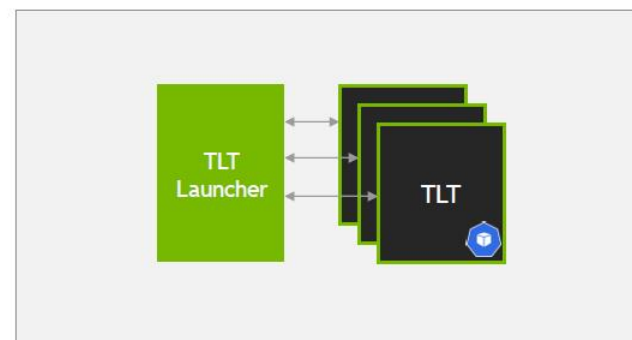
15 DOMAIN-SPECIFIC PRE-TRAINED MODELS
Accelerate time to market



MULTI GPU & MULTI NODE TRAINING
Speedup training time by 8x



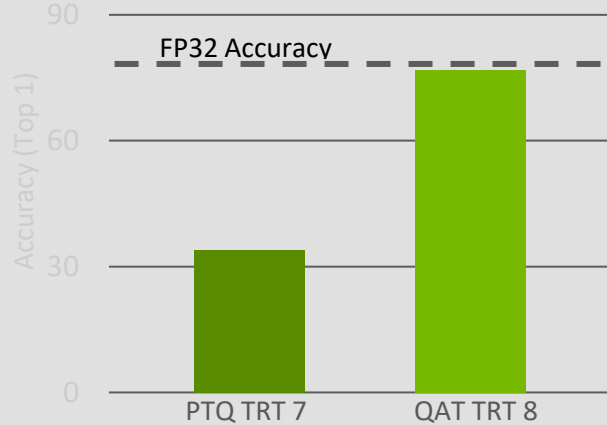
NEW MODEL ARCHITECTURES
EfficientNet, YOLOV4, Facial Landmark Detection, OCR



EASE OF USE
A Unified TLT launcher to manage and orchestrate your workflow

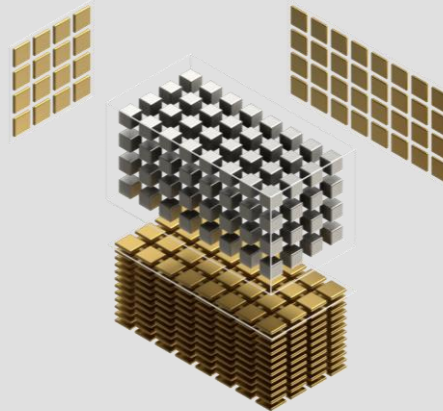
QAT, Sparsity | CNNs, RNNs, MLPs, Transformers

FP32 Accuracy with QAT INT8
(EfficientNet-B0)



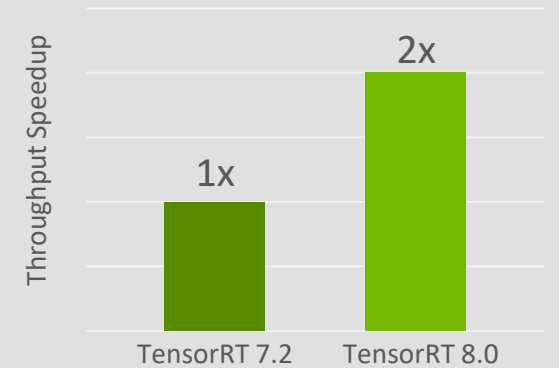
Quantization Aware Training delivers FP32 accuracy with INT8 precision

Upto 1.5x Faster inference on
Ampere GPUs



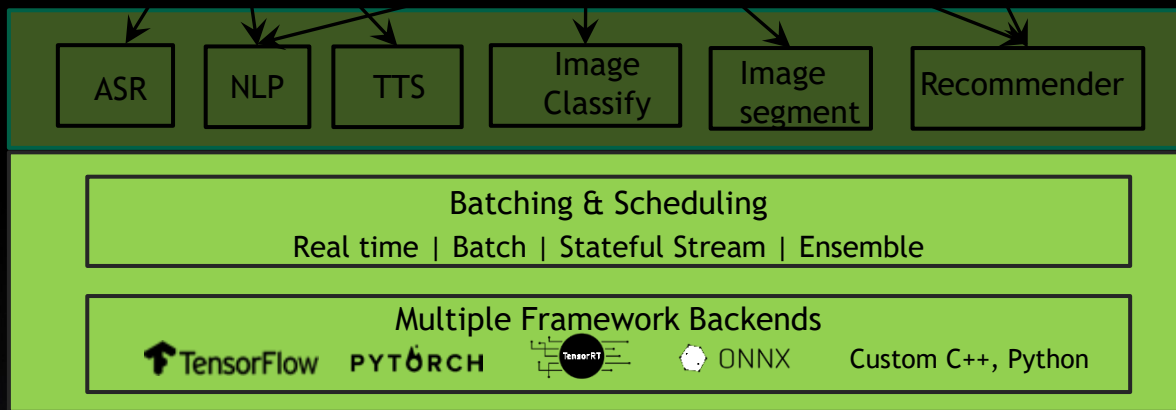
New Sparsity support for improved performance on Tensor Cores with INT8

Up To 2x Out-of-Box Performance
for Transformer Models



New optimizations for Transformer based models like BERT

To learn more, signup for the NVIDIA Developer Program from developer.nvidia.com/tensorrt

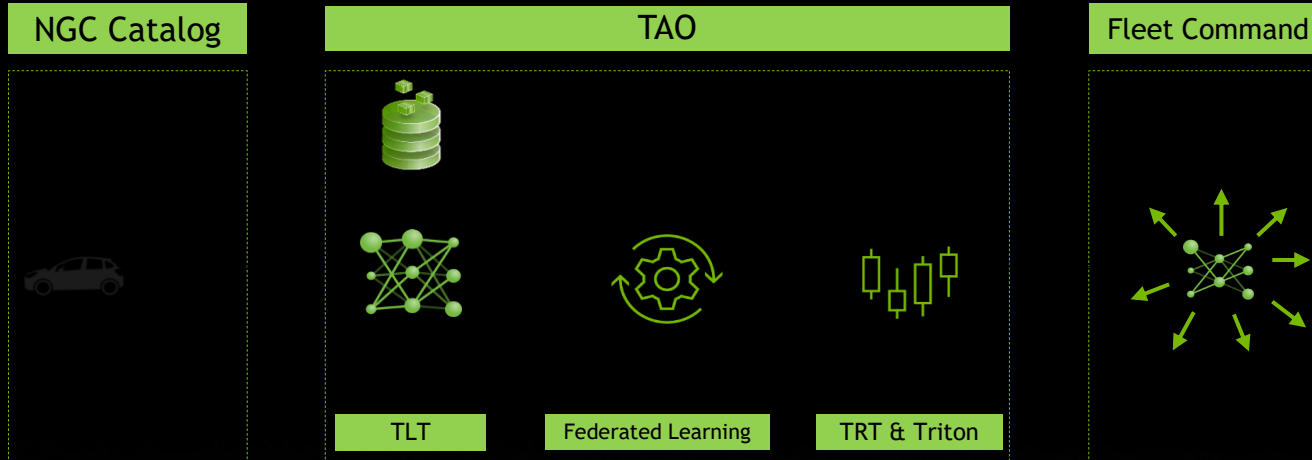


X86 CPU T4 GPU A100 GPU A100 MIG V100 GPU ARM CPU



Optimized for All Processors

Train | Adapt | Optimize



TRAIN

- Simplified framework simplifies AI development
- Train specific models in hours v. months

ADAPT & OPTIMIZE

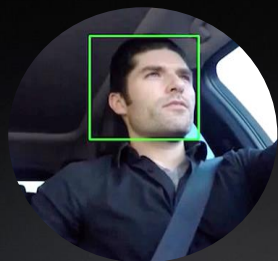
- Increase model accuracy with federated learning
- Optimize with TensorRT

DEPLOY

- Deploy from anywhere to anywhere
- Effortless management and security

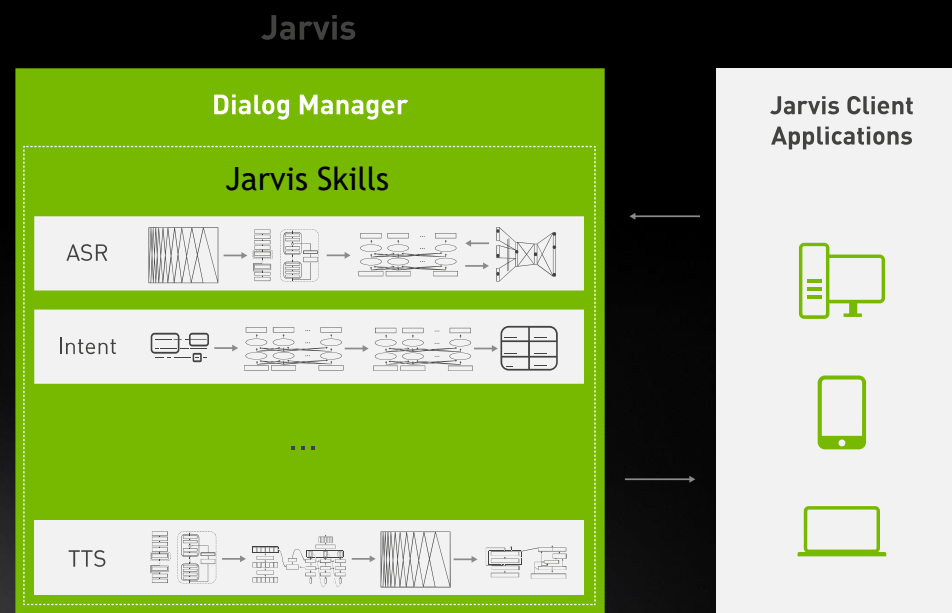
NVIDIA TAO availability: 2H, 2021

CONVERSATIONAL AI-JARVIS



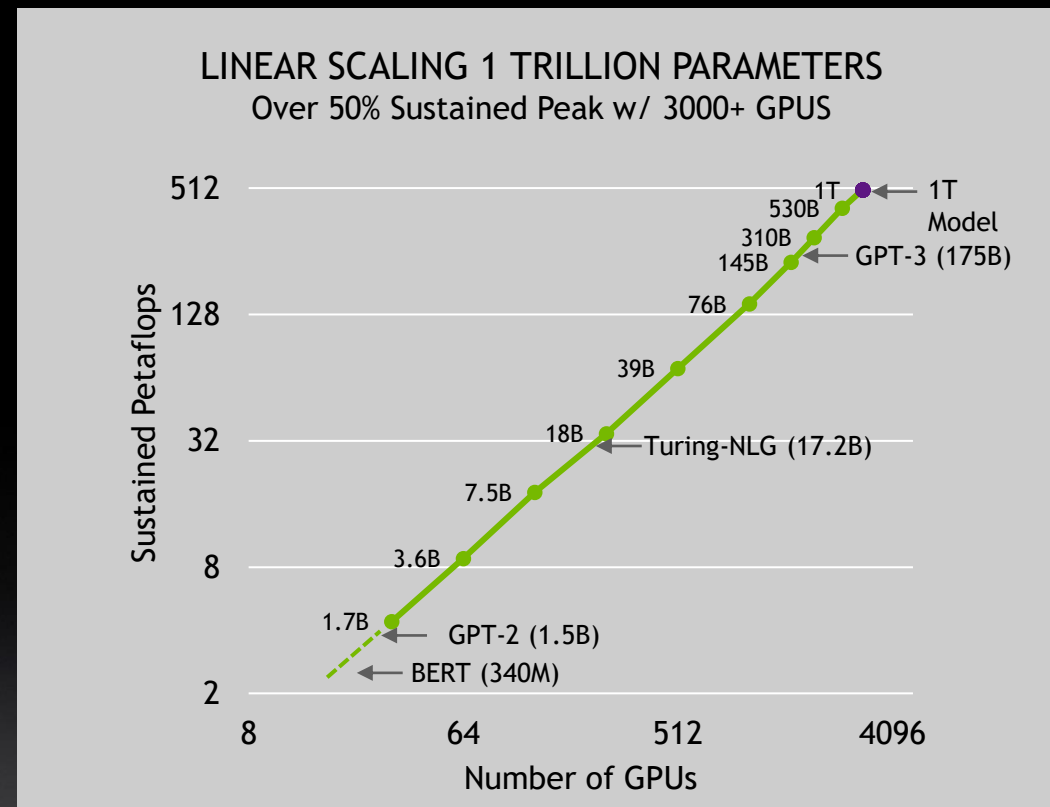
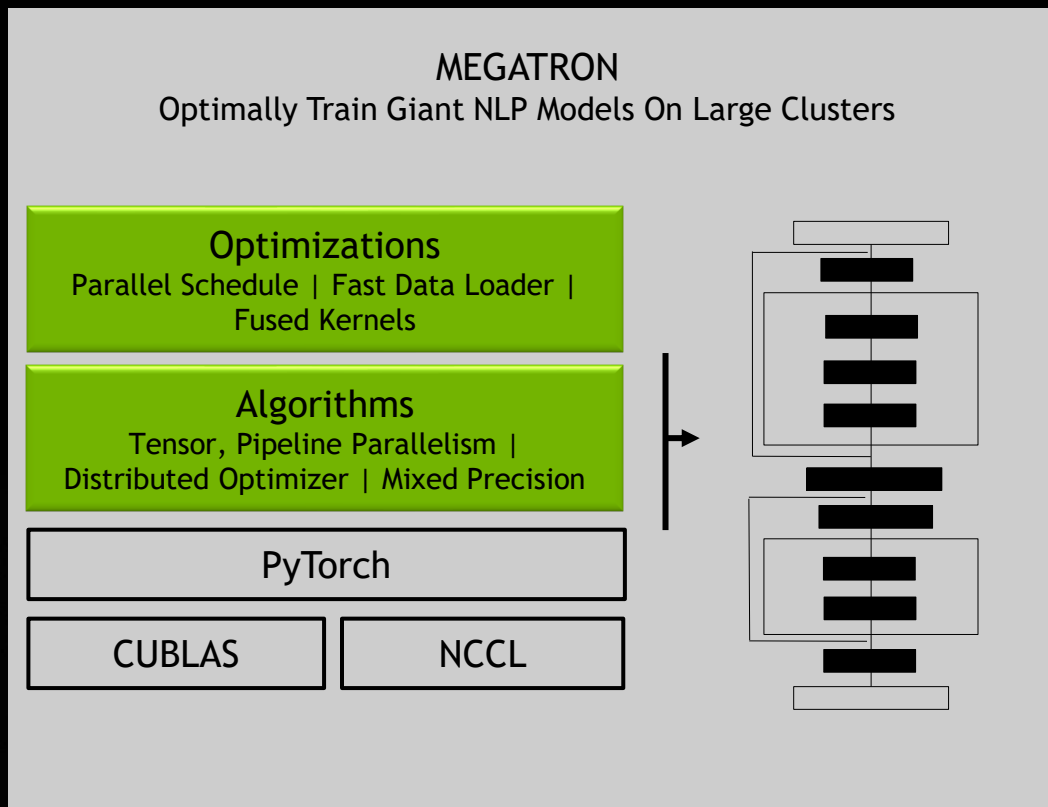
RETAIL ASSISTANTS
12M Retail Stores

IN-CAR ASSISTANTS
75M New Cars per Year



NVIDIA MEGATRON TRAINING FRAMEWORK FOR GIANT TRANSFORMERS

Train Multi-Trillion Parameter Models

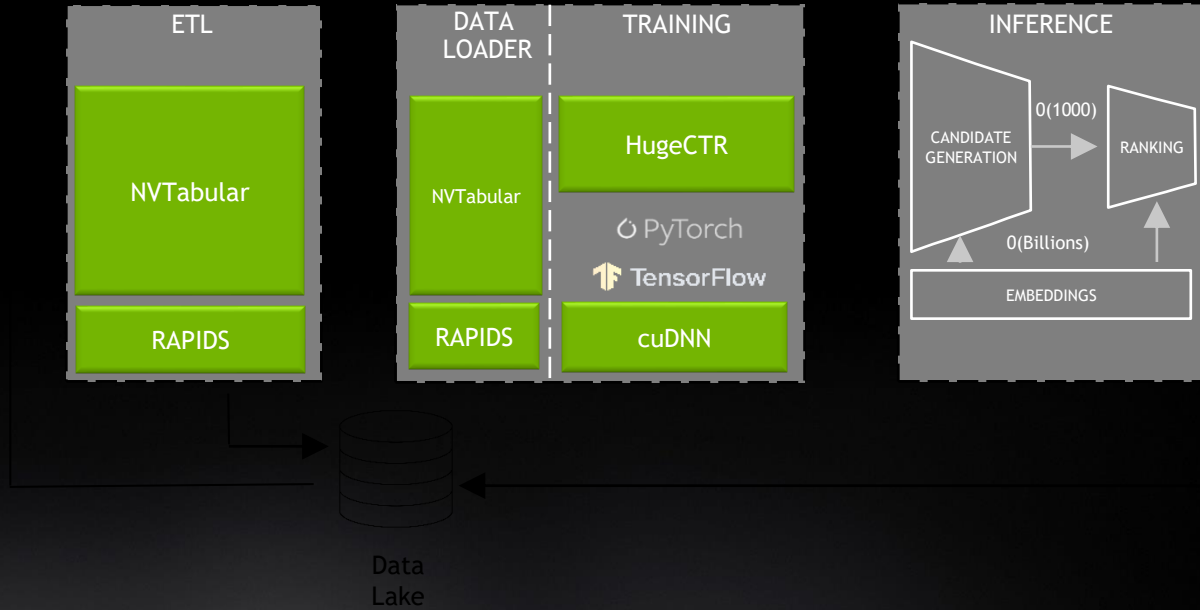


NVIDIA MERLIN END-TO-END ACCELERATED RECOMMENDER SYSTEM

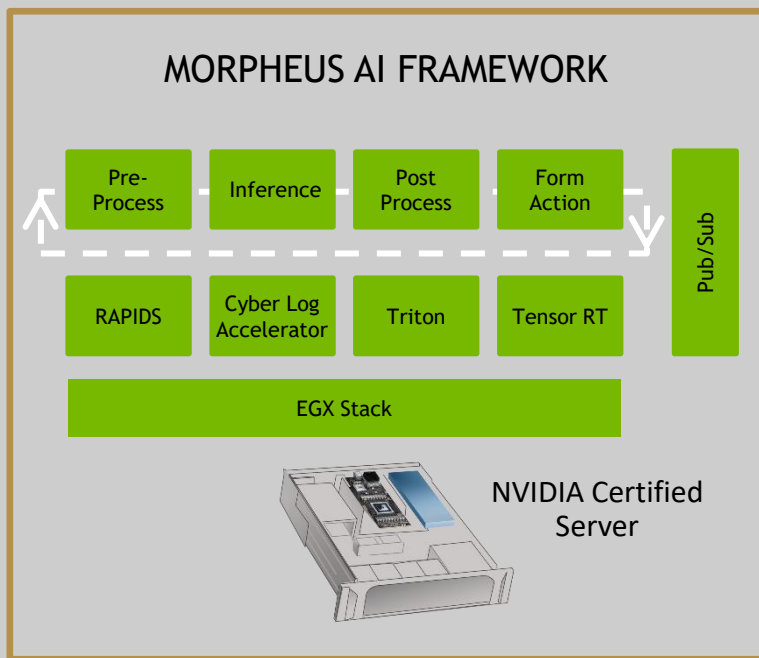


MUSIC
\$1.5B Purchased 2020

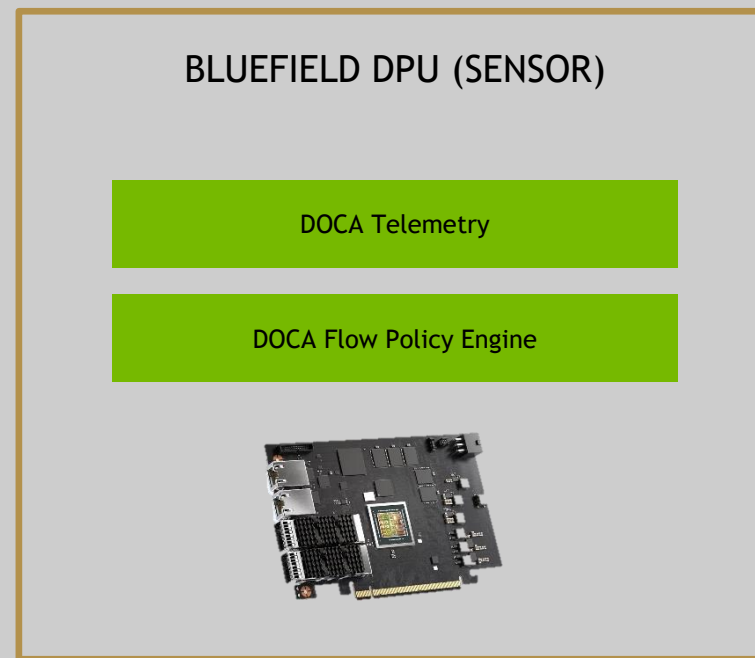
VIDEO
1B Streaming Subscriptions



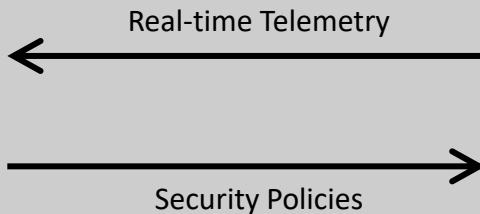
Open AI Compute Framework to Accelerate Cyber Ecosystem



ORCHESTRATE
Automated, Real-Time Threat
Detection & Response



SENSE
Penalty-Free Visibility &
Micro-Segmented Security



A PLATFORM BUILT FOR THE FUTURE

Multi-layer for multiple audiences

- ▶ Cloud-native
- ▶ Multi-GPU Enabled
- ▶ Open standards for cross-team, tool and workflow collaboration, built on Pixar's USD
- ▶ Scalable computing to address all workloads
- ▶ Works on all NVIDIA RTX™ solutions, from laptops to data centers



ACCELERATE SCIENTIFIC DISCOVERIES WITH NVIDIA OMNIVERSE

ParaView Omniverse Connector

Import Assets from Various Sources

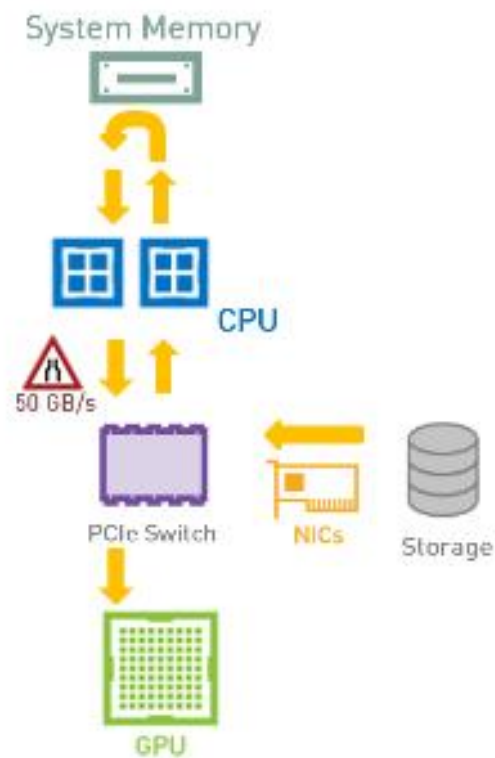
Real-time Collaboration Between Remote Teams

Photorealistic Scientific Visuals

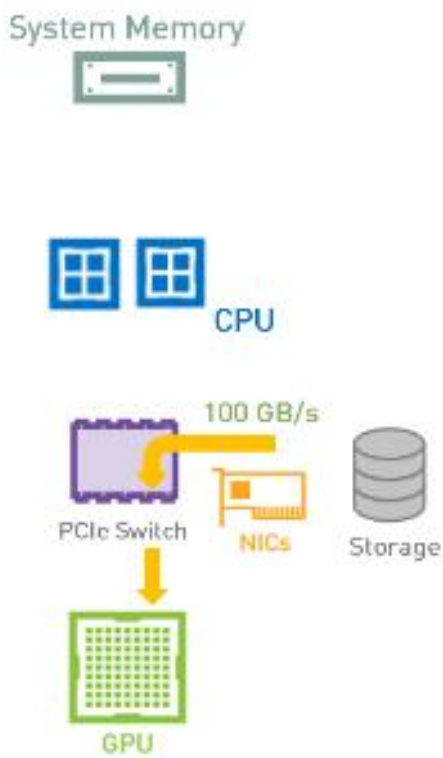
Available Now in Open Beta



WHAT IS GPUDIRECT STORAGE?

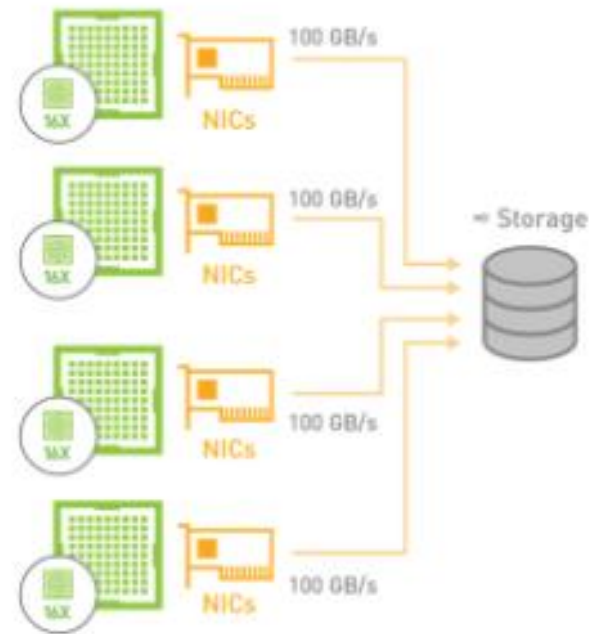


Without GPUDirect Storage



With GPUDirect Storage

- Higher Bandwidth
- Lower Latency

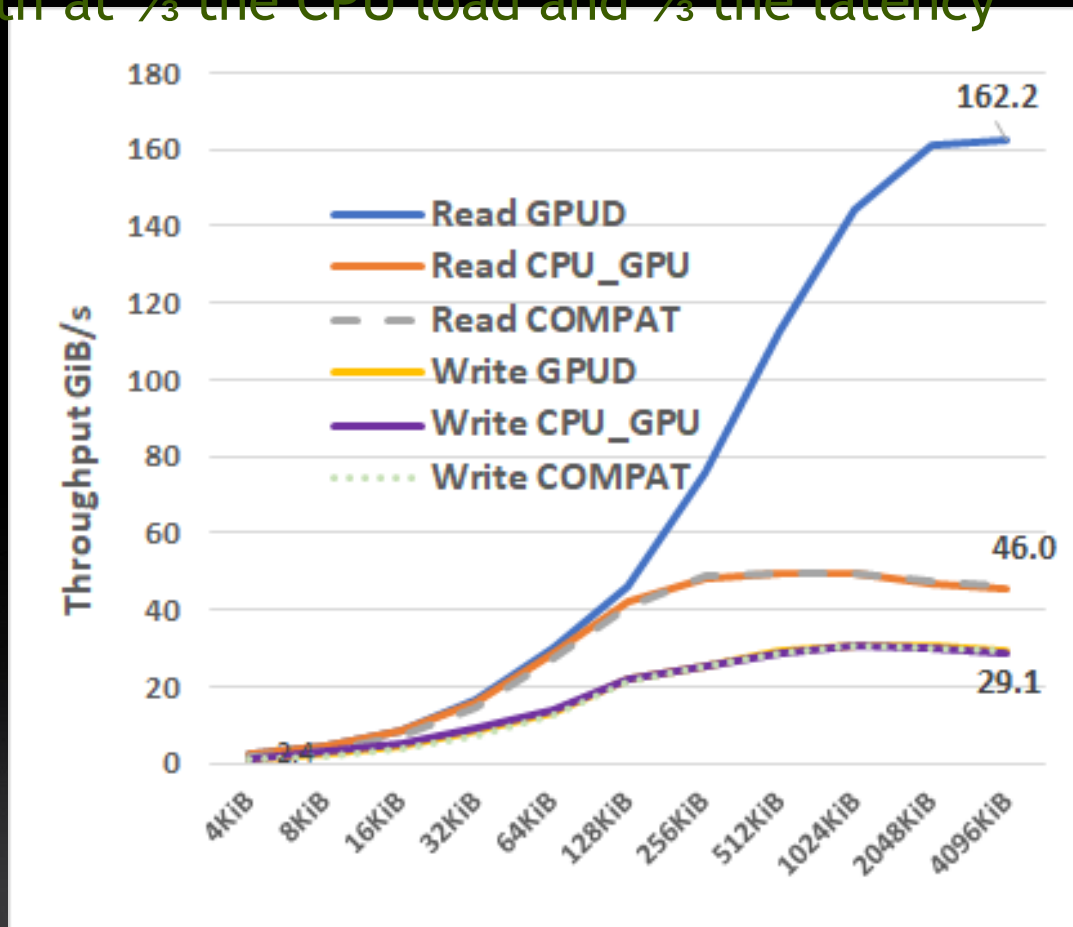


Scaling up with GPUDirect Storage



GDS provides 3x the bandwidth at 1/3 the CPU load and 1/3 the latency

- GDS path
- CPU's data movement
- GDS reads 162 GiB/s @ 14% CPU, 19ms
- CPU reads 50 GiB/s @ 49% CPU, 66ms
- GDS writes 30 GiB/s @ 3% CPU, 12ms
- CPU writes 29 GiB/s @ 4% CPU, 52ms



GA



BETA



Emerging

Early Exploration



NVIDIA DOCA

Enabling Broad DPU Partner Ecosystem

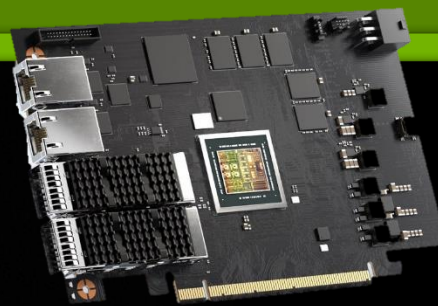
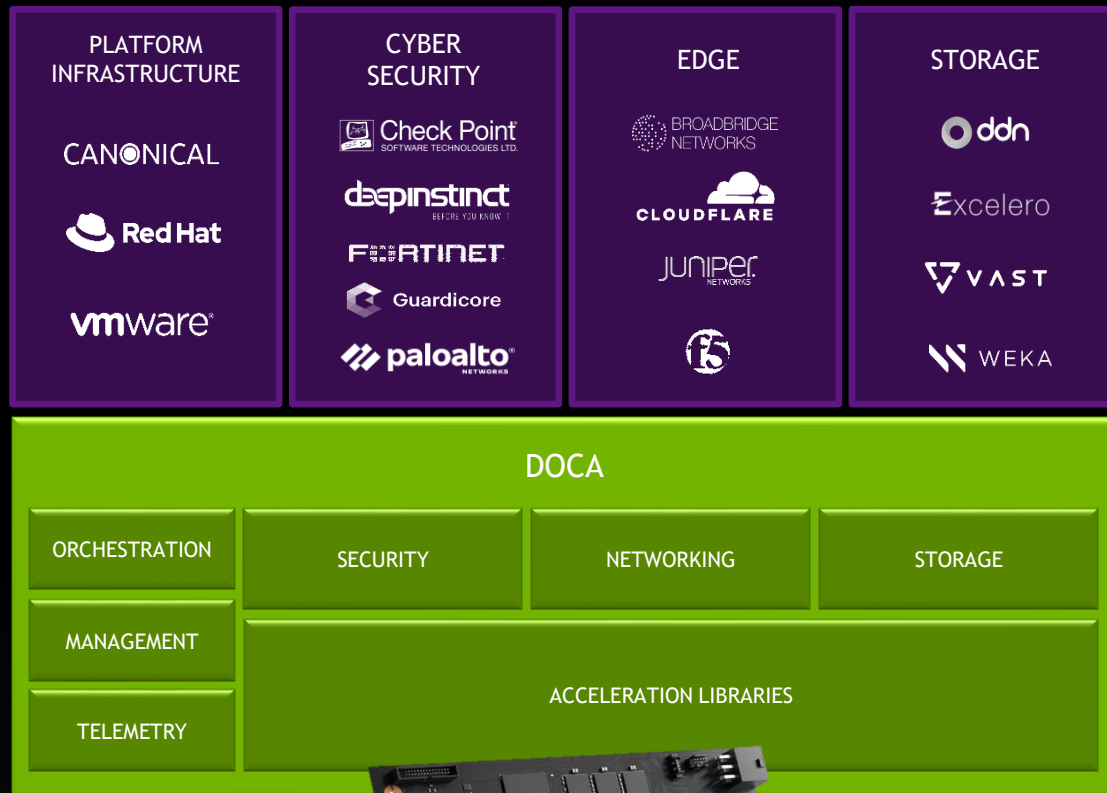
Software application framework for BlueField DPUs

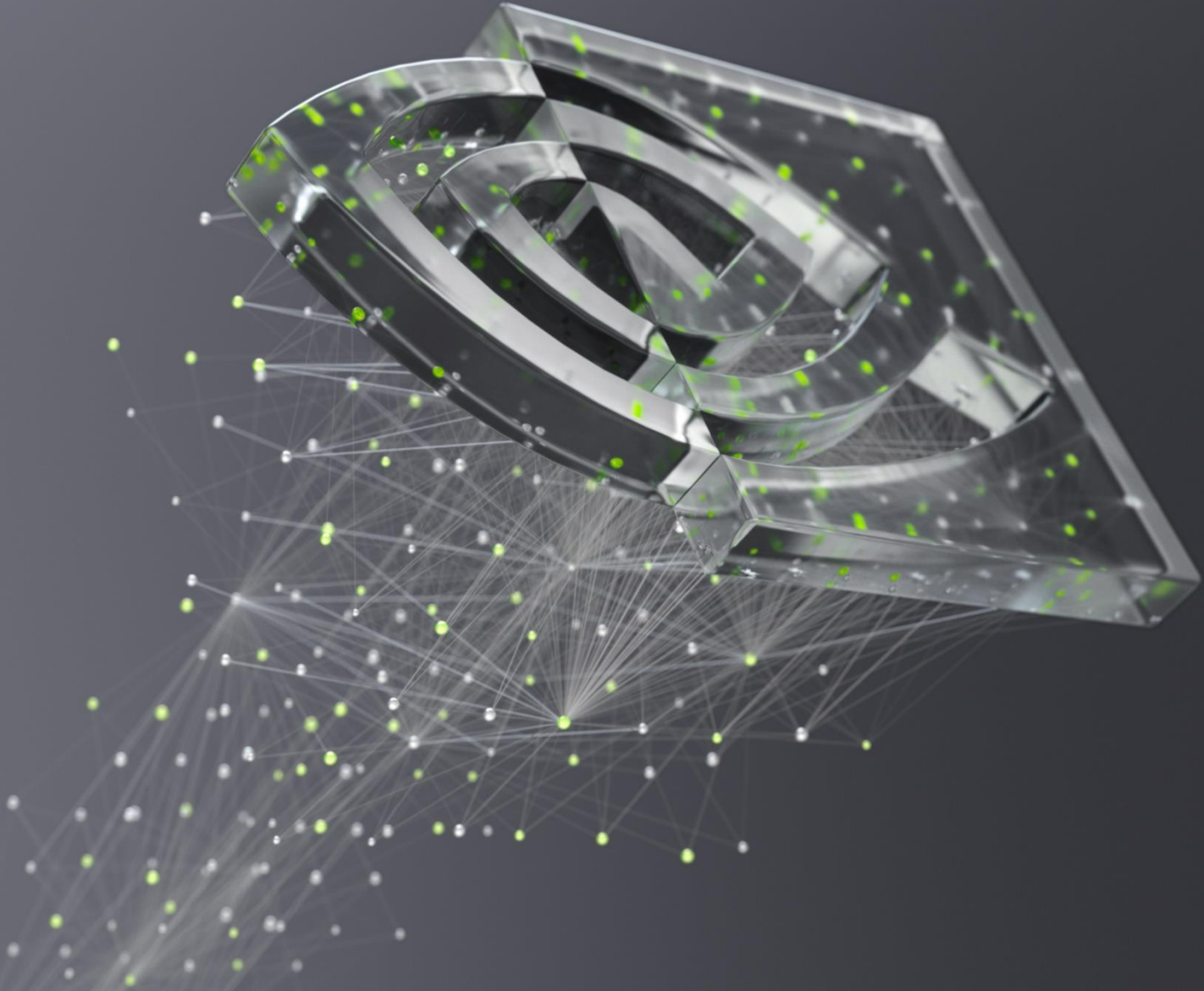
DOCA is for DPUs what CUDA is for GPUs

Protects developer investment for future DPUs

Certified reference applications, APIs & partner solutions

Rich partner ecosystem across industries and workloads





nVIDIA