EFFICIENT AT SCALE PERFORMANCE USING NVIDIA ACCELERATED SOLUTIONS FOR HPC AND AI 11

DR. GABRIEL NOAJE PRINCIPAL SOLUTIONS ARCHITECT, APAC SOUTH, NVIDIA **GNOAJE@NVIDIA.COM** JULY 26TH, 2022







GRAND CHALLENGES OF THE 21ST CENTURY





MILLION-X SPEEDUP FOR INNOVATION AND DISCOVERY Combination of Accelerated Computing, Data Center Scale and Al



SIMULATION + AI: MILLION-X SPEEDUP FOR INNOVATION AND DISCOVERY Combination of Accelerated Computing, Data Center Scale and AI

WORKLOADS OF THE MODERN SUPERCOMPUTER

EDGE

HPC + AI

SIMULATION

DIGITAL TWIN

QUANTUM COMPUTING

TRAINING NEURAL NETWORKS USING BOTH DATA AND THE GOVERNING EQUATIONS

NVIDIA MODULUS Physics Machine Learning Platform

Get started today with <u>NVIDIA Modulus</u>

INDUSTRIAL HPC NETL: 10,000X Faster Build Of highfidelity surrogate models

Kinetic Vision: Design Optimization Using

ACCELERATE QUANTUM COMPUTING RESEARCH WITH NVIDIA CUQUANTUM Research the Computer of Tomorrow with the Most Powerful Computer Today

COMPUTING OF TOMORROW

Seamless Acceleration Tensor Cores, Enhanced L2\$ & SMEM

NVIDIA PERFORMANCE LIBRARIES Major Directions

Performance: cuFFTMp vs. State-of-the-Art on Summit

MULTI-NODE MATH LIBRARIES cuFFTMp: Fast Fourier Transforms at scale

cuFFTMp State-of-the-Art

cuFFTMp

A distributed-memory multi-node and multiGPU solution for solving FFTs at scale.

EA release available in Fall '21 https://developer.nvidia.com/cudamathlibraryea

Initial release to 2D & 3D with Slab composition

FLEXIBLE HPC PROGRAMMING OPTIONS Develop How You Want, Where You Want

23X MORE PERFORMANCE IN 3.5 YEARS NVIDIA AI Delivers Continuous Gains With SW and At-Scale Improvements

Results normalized for throughput due to higher accuracy requirements on latest round of MLPerf Training 2.0. MLPerf ID 0.5/0.7/2.0 comparison: BERT: 0.7-56/0.7-38/2.0-2106 | ResNet50v1.5: 0.5-17/0.7-55/0.7-37/2.0-2107 | Mask R-CNN: 0.5-22/0.7-48/0.7-28/2.0-2099 MLPerf name and logo are trademarks. See <u>www..org</u> for more information.

Middle Panel: Ansys Fluent 2022 beta 1 running 105M cell car model | CPU baseline Intel Ice Lake 8380 with 80 cores | GPU A100 A100 PCIe 80GB | Right Panel: Siemens Simcenter STAR-CCM+ 2022.1 LeMans 104M model | CPU baseline AMD 7742 Rome |. GPU V100 PCIe 16 GB | A100 A100 PCIe 80GB

ACCELERATING INDUSTRIAL HPC SIMULATIONS

HIGHEST AI AND HPC PERFORMANCE

4PF FP8 (6X) | 2PF FP16 (3X) | 1PF TF32 (3X) | 60TF FP64 (3X) 3TB/s (1.5X), 80GB HBM3 memory

TRANSFORMER MODEL OPTIMIZATIONS

6X faster on largest transformer models

HIGHEST UTILIZATION EFFICIENCY AND SECURITY

7 Fully isolated & secured instances, guaranteed QoS 2nd Gen MIG | Confidential Computing

FASTEST, SCALABLE INTERCONNECT

900 GB/s GPU-2-GPU connectivity (1.5X) up to 256 GPUs with NVLink Switch | 128GB/s PCI Gen5

NVIDIA H100

Unprecedented Performance, Scalability, and Security for Every Data Center

Unique Versatile Architecture Combines H100 GPU and ConnectX-7 SmartNIC on a single board

Better I/O Performance High-speed, contention-free GPU-Network data transfer

Optimized Design Inherently balanced scale-out architecture

Cost Savings High performance with mainstream servers

Application Ready Acceleration benefits without application modification

NVIDIA H100 CNX Unprecedented performance for GPU-powered, IO-intensive workloads

Multi-node Training

350W | 80GB | 400 Gb/s Eth or IB PCIe Gen 5 within board and to host 2-Slot FHFL | NVLink

5G vRAN

Al on 5G

HIGHEST ACCELERATED PERFORMANCE Grace CPU plus Hopper GPU Acceleration

~600GB MEMORY AVAILABLE TO GPU

Enables Giant AI Models for Training & Inference

HIGHEST MEMORY BANDWIDTH 3.5GB/s LPDDR5x and HBM3

NEW 900GB/S COHERENT INTERFACE

NVLink-C2C connecting Grace to Hopper

15X HIGHER SYSTEM MEMORY BANDWIDTH TO GPU NVLink-C2C vs PCle

RUNS FULL NVIDIA COMPUTING STACKS

RTX, HPC, AI, Omniverse

AVAILABLE 1H 2023

FP8, FP16, TF32 performance include sparsity. X-factor compared to A100

GRACE HOPPER SUPERCHIP Built for Giant Scale AI and HPC

HIGHEST CPU PERFORMANCE Superchip Design with 144 high-performance Armv9 Cores Estimated Specrate2017_int_base of over 740

HIGHEST MEMORY BANDWIDTH

World's first LPDDR5x memory with ECC, 1TB/s Memory Bandwidth

HIGHEST ENERGY EFFICIENCY

2X Perf/Watt, CPU Cores + Memory in 500W

2X PACKING DENSITY

2x density of DIMM based designs

RUNS FULL NVIDIA COMPUTING STACKS RTX, HPC, AI, Omniverse

AVAILABLE 1H 2023

GRACE CPU SUPERCHIP The CPU for AI and HPC Infrastructure

DELIVERING UP TO 5X HPC HIGHER PERFORMANCE

GRACE HOPPER SUPERCHIP

Grace Hopper Superchip and Grace CPU Superchip Pre-silicon ESTIMATES ONLY: | Performance comparisons: (left) DGX-A100 (4x A100 SXM4 80 GB) | CP2K version 9.1 Random Phase Approximation dataset H2O-128-RI-dRPA-TZ.inp | ABINIT version 8.10.3 dataset Titanium 255 Atoms, LOBPCG algorithm | (Right) Traditional CPU Intel Ice Lake 8380 | BWA MEM2 v2.2.1 from the git repo running Full Human Genome (human_g1k_v37) | OpenFOAM version 2112 for model Large Motorbike v1912

GRACE CPU SUPERCHIP

GRACE REFERENCE DESIGNS FOR MODERN DATA CENTER WORKLOADS

Grace Hopper Superchip BlueField-3 OEM Defined IO / 4th Gen NVLink

DGX SUPERPOD WITH NVLINK SWITCH SYSTEM Purpose Built for Giant NLP, DLRM, Scientific Computing (3DFFT), MoE Models

Availability

2023 Only sold as part of DGX SuperPOD solution for enterprises and CSPs.

Scale Up AI and HPC

Delivers massive collective communication performance 9X bandwidth compared to HDR InfiniBand

H100 CLUSTER (1 SCALABLE UNIT)

57.6 TB/s Bi-section Bandwidth 32 servers | 18 NVLink switches | 1,152 NVLink optical cables

Scale Out Compute and Storage Communications

uses Quantum-2 NDR InfiniBand

QUANTUM-2 INFINIBAND SWITCH

Stores et al 1 Stores et al

Cloud Native Supercomputing Platform SHARP In-Network Computing Higher Scalability

CONNECTX-7 SMARTNIC

Intelligent Offloads Precision Timing Software Defined Networking

BLUEFIELD-3 / -X DPU

Intelligent Offloads Precision Timing Software Defined Networking

NVIDIA CLOUD NATIVE SUPERCOMPUTING

In-Network Computing

Zero Trust Security

Computational Storage

Enhanced Telemetry

SKYWAY GATEWAY InfiniBand to Ethernet Low Latency Load Balancing UFM

Monitoring, Management, Orchestration Predictive Maintenance Anomaly Detection

RADICAL IMPROVEMENTS IN CPU EFFICIENCY

Node Configuration	2x 96c Genoa Estimated	Single Grace Su Projected
General-Purpose FP64	7.4 Gflops	7.4 Gflop
Memory Bandwidth	1 TB/s	1 TB/s
CPU + Memory Power	900W	500W
Bandwidth per Core	5.3 GB/s	7.0 GB/

For 5 Petaflops of HPC capacity (~2 GB/core):

Datacenter Power Needed

5-Year TCO Savings

(\$.25/KWH, Estimated: Power, Cooling, CapEx)

Number of Cabinets

(300 KW LC or similar: 2.5 of Genoa, 1.5 of Grace)

Lifetime CO₂ Savings

(Assuming 5-Year Deployment)

NVIDIA CONFIDENTIAL. DO NOT DISTRIBUTE.

Grace Delivers On Both Compute And Memory Bandwidth

1 x	1.6x
\$2.3 Milli	
2	3
10,000 Metri	

Grace Superchip vs. 2S X86-64 + Memory

1.8x Perf / Watt.

Big Results

35% Lower OpEx

50% Improved Density

Save Millions in TCO

...and CO₂ emissions equivalent to 87,000 trees.

LIQUID-COOLED A100 TENSOR CORE GPU PORTFOLIO

A100 Offered with Liquid-Cooled Options

A100 HGX Shipping from OEMs

ENERGY & SPACE EFFICIENCY A100 PCIE LIQUID COOLED VS AIR COOLED

A100 PCIe Q3 2022

<1.2 PUE

EASILY MEET EFFICIENCY TARGETS

NVIDIA NEXT-GEN COMPUTING PLATFORM POWERING THE NEXT WAVE OF AI SUPERCOMPUTERS

Hopper + X86 systems: University of Tsukuba, Bristol, and TACC Grace Hopper/Grace CPU Superchips systems: CSCS and LANL

Quantum 200Gb/s InfiniBand

3 CHIPS. YEARLY LEAPS. ONE ARCHITECTURE.

