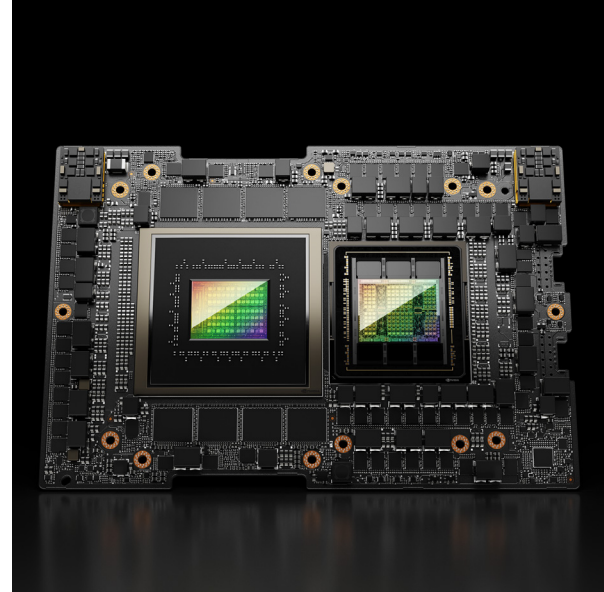




NVIDIA GH200 Grace Hopper Superchip

The breakthrough processor for large-scale AI and high-performance computing (HPC) applications.



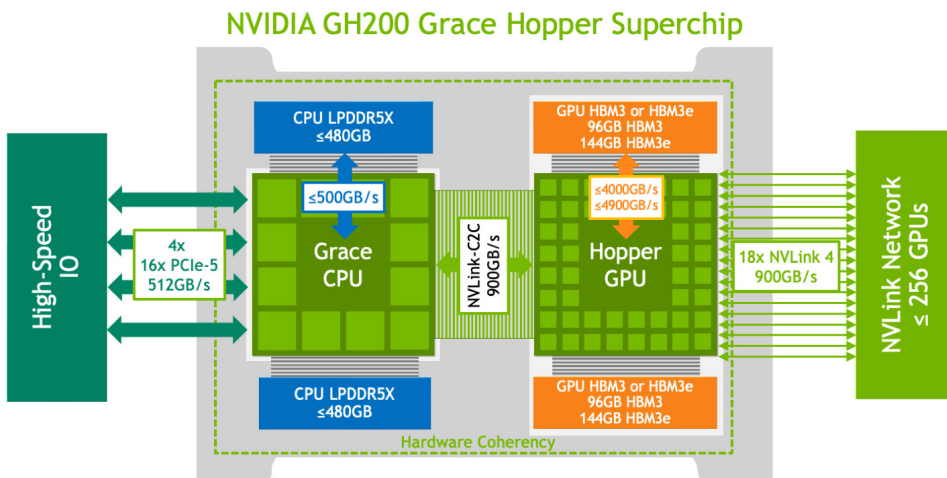
The World's Most Versatile Computing Platform

The NVIDIA Grace Hopper™ architecture brings together the groundbreaking performance of the NVIDIA Hopper™ GPU with the versatility of the NVIDIA Grace™ CPU in a single superchip, connected with the high-bandwidth, memory-coherent NVIDIA® NVLink® Chip-2-Chip (C2C) interconnect.

NVIDIA NVLink-C2C is a memory-coherent, high-bandwidth, and low-latency interconnect for superchips. The heart of the GH200 Grace Hopper Superchip, it delivers up to 900 gigabytes per second (GB/s) of total bandwidth, which is 7X higher than PCIe Gen5 lanes commonly used in accelerated systems. NVLink-C2C enables applications to oversubscribe the GPU's memory and directly utilize the Grace CPU's memory at high bandwidth. With up to 480GB of LPDDR5X CPU memory per GH200 Grace Hopper Superchip, the GPU has direct access to 7X more fast memory than HMB3 or almost 8X more fast memory than HBM3e, depending on the GH200 memory configured. GH200 can be easily deployed in standard servers to run a variety of inference, data analytics, and other compute- and memory-intensive workloads. GH200 can also be combined with the NVIDIA NVLink Switch System, with all GPU threads running on up to 256 NVLink-connected GPUs and able to access up to 144 terabytes (TB) of memory at high bandwidth.

Key Features

- > 72-core NVIDIA Grace CPU
- > NVIDIA H100 Tensor Core GPU
- > Up to 480GB of LPDDR5X memory with error-correction code (ECC)
- > Supports 96GB of HBM3 or 144GB of HBM3e
- > Up to 624GB of fast-access memory
- > NVLink-C2C: 900GB/s of coherent memory



Power and Efficiency With the Grace CPU

The NVIDIA Grace CPU delivers 2X the performance per watt of conventional x86-64 platforms and is the world's fastest Arm® data center CPU. The Grace CPU was designed for high single-threaded performance, high-memory bandwidth, and outstanding data-movement capabilities. The NVIDIA Grace CPU combines 72 Neoverse V2 Armv9 cores with up to 480GB of server-class LPDDR5X memory with ECC. This design strikes the optimal balance of bandwidth, energy efficiency, capacity, and cost. Compared to an eight-channel DDR5 design, the Grace CPU LPDDR5X memory subsystem provides up to 53 percent more bandwidth at one-eighth the power per gigabyte per second.

Performance and Speed With the Hopper H100 GPU

The H100 Tensor Core GPU is NVIDIA's ninth-generation data center GPU, and it delivers an order-of-magnitude performance leap for large-scale AI and HPC over the prior-generation NVIDIA A100 Tensor Core GPU. The NVIDIA H100 based on the new Hopper GPU architecture features multiple innovations:

- New fourth-generation Tensor Cores perform faster matrix computations than ever before on an even broader array of AI and HPC tasks.
- A new Transformer Engine enables H100 to deliver up to 9X faster AI training and up to 30X faster AI inference compared to the prior GPU generation.
- Secure Multi-Instance GPU (MIG) partitions the GPU into isolated, right-size instances to maximize quality of service (QoS) for smaller workloads.

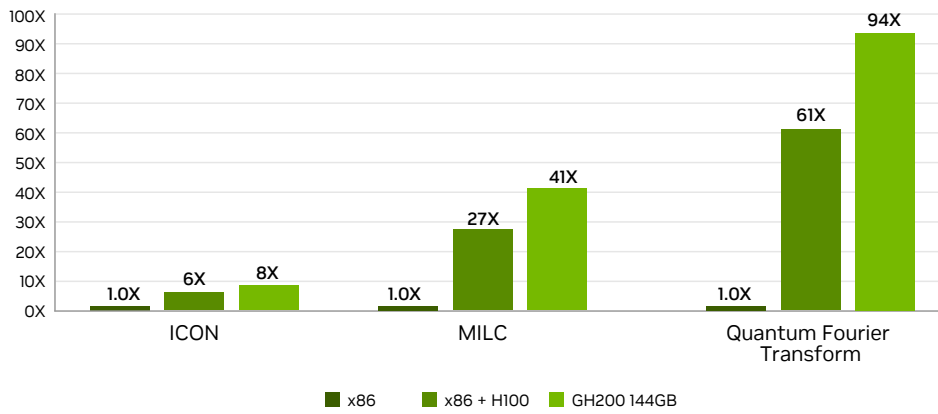
The Power of Coherent Memory

NVLink-C2C memory coherency increases developer productivity, performance, and the amount of GPU-accessible memory. CPU and GPU threads can concurrently and transparently access both CPU and GPU resident memory, allowing developers to focus on algorithms instead of explicit memory management. Memory coherency lets developers only transfer the data they need and not migrate entire pages to and from the GPU. It also provides lightweight synchronization primitives across GPU and CPU threads by enabling native atomics from both the CPU and GPU. Fourth-generation NVLink allows accessing peer memory with direct loads, stores, and atomic operations, so accelerated applications can solve larger problems more easily than ever.

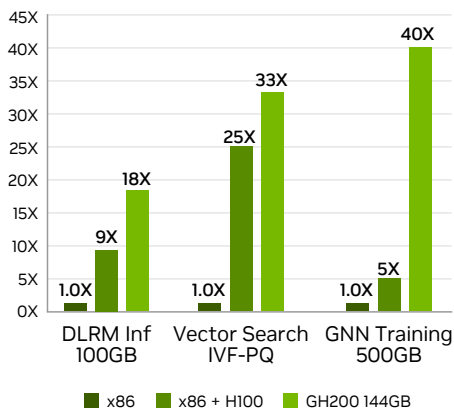
Class-Leading Performance for HPC and AI Workloads

The GH200 Grace Hopper Superchip is the first true heterogeneous accelerated platform for HPC and AI workloads. It accelerates any application with the strengths of both GPUs and CPUs while providing the simplest and most productive heterogeneous programming model to date, enabling scientists and engineers to focus on solving the world's most important problems. For AI inference workloads, GH200 Grace Hopper Superchips combine with NVIDIA networking technologies to provide the best TCO for scale-out solutions, letting customers take on larger datasets, more complex models, and new workloads using up to 624GB of fast-access memory. The NVIDIA GH200 also comes in a dual-GH200 configuration with two Grace Hopper Superchips fully connected by NVLink to deliver 288GB of HBM3e and 1.2TB of fast memory for both compute- and memory-intensive workloads.

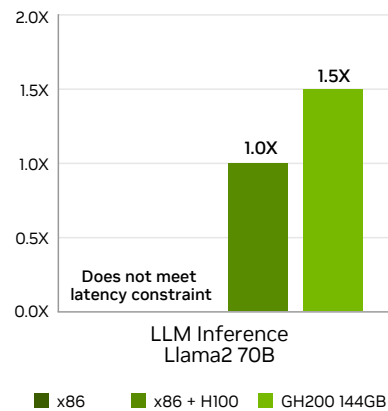
GH200 HPC Performance



GH200 AI Performance



GH200 LLM Performance



* Comparison: 2S Xeon Platinum 8480+, Xeon Platinum 8480+ and H100 Tensor Core GPU, and NVIDIA GH200 144GB: ICON (QUBICC 191 levels 80km radiation), MILC (APEX Medium), Quantum Fourier Transform, DLRM 100GB inference, vector search (batch size = 10,000 | queries = 10,000 over 85M vectors IVF-PQ) GNN (GraphSage OGB-100M papers dataset), Llama2 70B (batch size = 64 (GH200), 96 (2x H100s) | precision = FP8 | TensorRT-LLM - throughput per GPU).

Results subject to change.

Full NVIDIA Platform Support

The NVIDIA GH200 Grace Hopper Superchip extends the existing large and diverse ecosystem of 64-bit Arm processors. The very same containers, application binaries, and operating systems that run on other Arm products run on Grace Hopper without modification—only faster. And for customers who wish to leverage and build upon NVIDIA's software expertise, the NVIDIA Grace Hopper Superchip is supported by the full NVIDIA software stack, including the NVIDIA HPC, NVIDIA AI, and NVIDIA Omniverse™ platforms.

Product Specifications

Grace CPU	Feature
CPU core count	72 Arm Neoverse V2 cores
L1 cache	64KB i-cache + 64KB d-cache
L2 cache	1MB per core
L3 cache	114MB
Base frequency all-core single instruction, multiple data (SIMD) frequency	3.1GHz 3.0GHz
LPDDR5X size	Up to 480GB
Memory bandwidth	Up to 512GB/s
PCIe links	Up to 4x PCIe x16 (Gen5)

Hopper H100 GPU	Feature
FP64	34 teraFLOPS
FP64 Tensor Core	67 teraFLOPS
FP32	67 teraFLOPS
TF32 Tensor Core	989 teraFLOPS* 494 teraFLOPS
BFLOAT16 Tensor Core	1,979 teraFLOPS* 990 teraFLOPS
FP16 Tensor Core	1,979 teraFLOPS* 990 teraFLOPS
FP8 Tensor Core	3,958 teraFLOPS* 1,979 teraFLOPS
INT8 Tensor Core	3,958 TOPS* 1,979 TOPS
High-bandwidth memory (HBM) size	Up to 96GB 144GB HBM3e
Memory bandwidth	Up to 4TB/s Up to 4.9TB/s HBM3e
NVIDIA NVLink-C2C CPU-to-GPU bandwidth	900 GB/s bidirectional
Module thermal design power (TDP)	Programmable from 450W to 1000W (CPU + GPU + memory)
Form factor	Superchip module
Thermal solution	Air cooled or liquid cooled

* With sparsity

Ready to Get Started?

To learn more about the NVIDIA Grace Hopper Superchip, visit [nvidia.com/grace-hopper-superchip/](https://www.nvidia.com/grace-hopper-superchip/)

To download the Grace Hopper architecture whitepaper, visit resources.nvidia.com/en-us-grace-cpu/nvidia-grace-hopper

© 2024 NVIDIA Corporation and affiliates. All rights reserved. NVIDIA, the NVIDIA logo, ConnectX, DGX, HGX, Hopper, NGC, NVIDIA-Certified Systems, and NVLink are trademarks and/or registered trademarks of NVIDIA Corporation and affiliates in the U.S. and other countries. Other company and product names may be trademarks of the respective owners with which they are associated. 3154100. MAR24.

