# Graphics Card NVIDIA 94 GB Memory 2700 MHz 699-21010-0210-XXX

Model: 699-21010-0210-XXX

Manufacturer: NVIDIA

Product Code: 14101

## R811,775

**An Order-of-Magnitude Leap for Accelerated Computing**

The NVIDIA H100 Tensor Core GPU delivers exceptional performance, scalability, and security for every workload. H100 uses breakthrough innovations based on the NVIDIA Hopper™ architecture to deliver industry-leading conversational AI, speeding up large language models (LLMs) by 30X. H100 also includes a dedicated Transformer Engine to solve trillion-parameter language models.

**Securely Accelerate Workloads From Enterprise to Exascale**

**Transformational AI Training**

H100 features fourth-generation Tensor Cores and a Transformer Engine with FP8 precision that provides up to 4X faster training over the prior generation for GPT-3 (175B) models. The combination of fourth-generation NVLink, which offers 900 gigabytes per second (GB/s) of GPU-to-GPU interconnect; NDR Quantum-2 InfiniBand networking, which accelerates communication by every

GPU across nodes; PCIe Gen5; and NVIDIA Magnum IO™ software delivers efficient scalability from small enterprise systems to massive, unified GPU clusters.

Deploying H100 GPUs at data center scale delivers outstanding performance and brings the next generation of exascale high-performance computing (HPC) and trillion-parameter AI within the reach of all researchers.

**Real-Time Deep Learning Inference**

AI solves a wide array of business challenges, using an equally wide array of neural networks. A great AI inference accelerator has to not only deliver the highest performance but also the versatility to accelerate these networks.

H100 extends NVIDIA's market-leading inference leadership with several advancements that accelerate inference by up to 30X and deliver the lowest latency. Fourth-generation Tensor Cores speed up all precisions, including FP64, TF32, FP32, FP16, INT8, and now FP8, to reduce memory usage and increase performance while still maintaining accuracy for LLMs.

**Exascale High-Performance Computing**

The NVIDIA data center platform consistently delivers performance gains beyond Moore's law. And H100's new breakthrough AI capabilities further amplify the power of HPC+AI to accelerate time to discovery for scientists and researchers working on solving the world's most important challenges.

H100 triples the floating-point operations per second (FLOPS) of double-precision Tensor Cores, delivering 60 teraflops of FP64 computing for HPC. AI-fused HPC applications can also leverage H100's TF32 precision to achieve one petaflop of throughput for single-precision matrix-multiply operations, with zero code changes.

H100 also features new DPX instructions that deliver 7X higher performance over A100 and 40X speedups over CPUs on dynamic programming algorithms such as Smith-Waterman for DNA sequence alignment and protein alignment for protein structure prediction.

**Accelerated Data Analytics**

Data analytics often consumes the majority of time in AI application development. Since large

datasets are scattered across multiple servers, scale-out solutions with commodity CPU-only servers get bogged down by a lack of scalable computing performance.

Accelerated servers with H100 deliver the compute power—along with 3 terabytes per second (TB/s) of memory bandwidth per GPU and scalability with NVLink and NVSwitch™—to tackle data analytics with high performance and scale to support massive datasets. Combined with NVIDIA Quantum-2 InfiniBand, Magnum IO software, GPU-accelerated Spark 3.0, and NVIDIA RAPIDS™, the NVIDIA data center platform is uniquely able to accelerate these huge workloads with higher performance and efficiency.

**Enterprise-Ready Utilization**

IT managers seek to maximize utilization (both peak and average) of compute resources in the data center. They often employ dynamic reconfiguration of compute to right-size resources for the workloads in use.

H100 with MIG lets infrastructure managers standardize their GPU-accelerated infrastructure while having the flexibility to provision GPU resources with greater granularity to securely provide developers the right amount of accelerated compute and optimize usage of all their GPU resources.

**Built-In Confidential Computing**

Traditional Confidential Computing solutions are CPU-based, which is too limited for compute-intensive workloads such as AI at scale. NVIDIA Confidential Computing is a built-in security feature of the NVIDIA Hopper architecture that made H100 the world's first accelerator with these capabilities. With NVIDIA Blackwell, the opportunity to exponentially increase performance while protecting the confidentiality and integrity of data and applications in use has the ability to unlock data insights like never before. Customers can now use a hardware-based trusted execution environment (TEE) that secures and isolates the entire workload in the most performant way.

**Exceptional Performance for Large-Scale AI and HPC**

The Hopper Tensor Core GPU will power the NVIDIA Grace Hopper CPU+GPU architecture, purpose-built for terabyte-scale accelerated computing and providing 10X higher performance on large-model AI and HPC. The NVIDIA Grace CPU leverages the flexibility of the Arm® architecture to create a CPU and server architecture designed from the ground up for accelerated computing. The Hopper GPU is

paired with the Grace CPU using NVIDIA's ultra-fast chip-to-chip interconnect, delivering 900GB/s of bandwidth, 7X faster than PCIe Gen5. This innovative design will deliver up to 30X higher aggregate system memory bandwidth to the GPU compared to today's fastest servers and up to 10X higher performance for applications running terabytes of data.

**Supercharge Large Language Model Inference With H100 NVL**

For LLMs up to 70 billion parameters (Llama 2 70B), the PCIe-based NVIDIA H100 NVL with NVLink bridge utilizes Transformer Engine, NVLink, and 188GB HBM3 memory to provide optimum performance and easy scaling across any data center, bringing LLMs to the mainstream. Servers equipped with H100 NVL GPUs increase Llama 2 70B performance up to 5X over NVIDIA A100 systems while maintaining low latency in power-constrained data center environments.

**Enterprise-Ready: AI Software Streamlines Development and Deployment**

NVIDIA H100 NVL comes with a five-year NVIDIA AI Enterprise subscription and simplifies the way you build an enterprise AI-ready platform. H100 accelerates AI development and deployment for production-ready generative AI solutions, including computer vision, speech AI, retrieval augmented generation (RAG), and more. NVIDIA AI Enterprise includes NVIDIA NIMTM, a set of easy-to-use microservices designed to speed up enterprise generative AI deployment. Together, deployments have enterprise-grade security, manageability, stability, and support. This results in performance-optimized AI solutions that deliver faster business value and actionable insights.

Summary: Nvidia H100 NVL. Graphics processor family: NVIDIA, Graphics processor: H100, Processor frequency: 1080 MHz. Discrete graphics card memory: 94 GB, Graphics card memory type: High Bandwidth Memory 3 (HBM3), Memory bus: 6016 bit, Memory clock speed: 2619 MHz. Interface type: PCI Express x16 5.0. Cooling type: Passive | Vendor Homepage: https://www.nvidia.com/en-us/data-center/h100/ | CUDA: Yes | Graphics processor: H100 | Graphics processor family: NVIDIA | Parallel processing technology support: NVLink | Processor boost clock speed: 1785 | Processor frequency: 1080 | Discrete graphics card memory: 94 | Graphics card memory type: High Bandwidth Memory 3 (HBM3) | Memory bandwidth (max): 3938 | Memory bus: 6016 | Memory clock speed: 2619 | Interface type: PCI Express x16 5.0 | Dual Link DVI: No | TV tuner integrated: No | Bracket height: Full-Height (FH) | Cooling type: Passive | Form factor: Full-Height/Full-Length (FH/FL) | Number of slots: 2 | Product colour: Beige | Power consumption (max): 400 | Supplementary power connectors: 1x 16-pin | Non-operating relative humidity (non-condensing): 5 - 95 | Operating relative humidity (H-H): 5 - 85 | Operating temperature (T-T): 0 - 50 | Storage temperature (T-T): -40 - 75 | Weight:

1214 | Model: 699-21010-0210-XXX | Units per Shipping Box: 1 | Unit Calculated Volume: 0.00594 | Product Net Weight: 1.214 | Unit Calculated Weight: 1.86