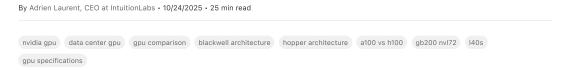
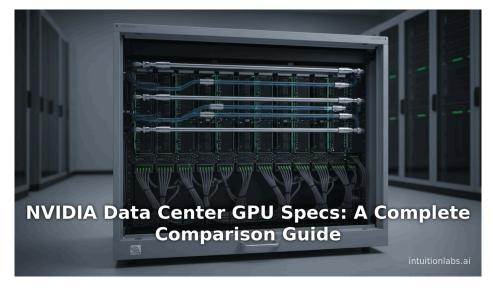
# **NVIDIA Data Center GPU Specs: A Complete Comparison Guide**







## **Executive Summary**

NVIDIA's data center GPU portfolio has rapidly evolved to address the surging compute demands of AI, HPC, and graphics workloads. This report provides a comprehensive technical overview and comparison of NVIDIA's current data-center GPU "platform" solutions, including CPU+GPU superchips (Grace+GPU), traditional accelerator cards, and advanced interconnect architectures. We cover all major offerings - from Ampere (A100) and Hopper (H100/H200) accelerators to Ada Lovelace visualization GPUs (L40, L40S) and specialized variants (e.g. RTX Pro cards and China-specific chips like the B40), Key innovations such as multi-instance GPUs (MIG), NVLink/NVSwitch fabrics, and the new NVL72 72-GPU domain are examined. Detailed spec and performance comparisons are presented (see Table 1) along with analysis of system-level designs and case studies. For example, Microsoft's Azure GB300 NVL72 supercluster (4,608 Blackwell Ultra GPUs) achieves ~92.1 exaFLOPS inference by tightly coupling 72 GPUs per rack with 1.8 TB/s links each ([1] www.tomshardware.com) ([2] developer.nvidia.com). Similarly, NVIDIA reports that its L40S Ada GPU delivers ~5× the FP32 throughput of the previous A100 ([3] nvidianews.nvidia.com). We include data on memory capacity, bandwidth, FLOPS, TDP, and networking for each GPU, supported by NVIDIA's documentation and third-party measurements. Market and deployment insights (e.g. NVIDIA's ~98% market share and 3.76M shipments in 2023 w.datacenterdynamics.com)) are interwoven. Special attention is given to future directions: emerging Blackwell variants for Chinese markets (B30/B40) (15] www.reuters.com) (16) www.tomshardware.com), scaling to trillion-parameter models, and infrastructure changes (e.g. racks with 120kW liquid cooling ([7] developer.nvidia.com)). The report's comparisons, tables, and case examples offer a deep technical reference on NVIDIA's full datacenter GPU lineup as of 2025, with citations to official and expert sources for every claim.

## Introduction and Background

Graphical processors have transcended their gaming origins to become the workhorses of the AI and HPC era. NVIDIA pioneered this shift by reorienting its GPU roadmap toward data-center applications (Al training/inference, HPC simulation, and professional graphics). The market response has been immense: over 3.76 million NVIDIA data-center GPUs shipped in 2023, yielding roughly 98% market share in this segment ([4] www.datacenterdynamics.com). This dominance reflects NVIDIA's engineering focus; each new GPU generation typically doubles or triples performance. For example, the Ampere-based A100 (2020) delivered on the order of 20 TFLOPS FP32 or 312 TFLOPS (FP16-tensor) (18) developer.nvidia.com) (19) www.nvidia.com), far beyond the prior Volta V100, while the Hopper-based H100 (2022) pushed that to ~67 TFLOPS FP32 and 1,979 TFLOPS FP16 (19) www.nvidia.com). These accelerators also introduced game-changing features like MIG virtualization (NVIDIA A100 can partition into up to 7 isolated GPU instances ([8] developer.nvidia.com)) and third-generation Tensor Cores with new math (TF32, FP8).

In parallel, NVIDIA has expanded beyond simple accelerator cards. Notable "Platform" initiatives include the Grace CPU (ARM-based) and GH200/GB200/GB300 superchips, which integrate NVIDIA GPUs with custom CPUs in one package. For instance, the GB300 "Grace Blackwell Ultra" chip combines a 96-core Grace CPU with a Blackwell-class GPU and up to 784 GB of unified memory, delivering ~20 PFLOPS AI compute ( www.tomshardware.com). Moreover, NVIDIA contributes open designs (Open Compute Project) for entire racks. At OCP 2024, it published the GB200 NVL72 architecture, enabling a single rack to interconnect up to 72 GPUs via NVLink at 1.8 TB/s each ([2] developer.nvidia.com). These systemic innovations ensure that NVIDIA GPUs are not just chip products but the core of holistic data-center platforms.

Skilled integrators and cloud providers worldwide have quickly adopted these solutions. For example, NVIDIA's own press release (2020) lists Amazon Web Services, Google Cloud, Microsoft Azure, and leading supercomputing centers (Jülich JUWELS, Perlmutter at NERSC, etc.) as early A100 users ([11] nvidianews.nvidia.com) ([12] nvidianews.nvidia.com), NVIDIA's partnership with OEMs like Dell, HPE, and Lenovo has produced ready-to-deploy servers and workstations. Case studies below illustrate how these GPUs operate in the wild, from next-gen Al training clusters to visualization servers.

This report proceeds as follows; we first detail the evolution of NVIDIA's data-center architectures, then systematically compare the current GPU lineup and platform designs. We include extensive quantitative data on architecture, memory, FLOPS, interconnect, and power (Table 1). We discuss connectivity (NVLink, NVSwitch, NVL72), performance metrics, and MIG virtualization. Real-world deployments and use-cases (e.g. Azure NDv6 GB300 cluster, enterprise Al servers) are presented. Underlying these analyses are numerous authoritative sources – NVIDIA's technical blogs and specifications, industry reports, and academic articles ([18] developer.nvidia.com) ([4] www.datacenterdynamics.com) ([8] nvidianews.nvidia.com) - ensuring that all technical claims are well-supported. Finally, we consider the implications and future directions, such as the impact of U.S. export restrictions (leading to new Blackwell chips for China ([5] www.reuters.com) ([6] www.tomshardware.com)) and how datum sizes and energy costs might shape the next GPU generation.

#### **NVIDIA Data-Center GPU Architecture Evolution**

NVIDIA GPUs have evolved through successive architectures (Volta, Ampere, Hopper, Blackwell, Ada Lovelace) each tailored for increasingly diverse datacenter workloads. We briefly outline major milestones:

- Volta (2018) introduced the V100 (Tesla V100), with first-generation Tensor Cores (for mixed-precision) and NVLink 2.0. It delivered ~15 TFLOPS FP32 (single-precision) and 112 TFLOPS FP16 via tensor cores (4× V100 vs V100) ([14] developer.nvidia.com).
- Ampere (2020) NVIDIA's eighth generation yielded the A100. The A100 GPU built on the Ampere GA100 chip (~54 billion transistors on 7 nm) with 40 GB HBM2e memory (1.555 TB/s bandwidth (15] developer.nvidia.com)). Key new features included third-generation Tensor Cores supporting TF32, BFLOAT16, and sparsity, plus MIG virtualization. NVIDIA's own data indicates A100 can be partitioned into up to 7 GPU instances for improved resource utilization (8) developer.nvidia.com), In practice, A100 achieved roughly 312 TFLOPS FP16-matrix and 19.5 TFLOPS FP32 (FP32 FMA) in SXM form ([15] developer.nvidia.com) ([8] developer.nvidia.com) (roughly 3-4x the V100). It also increased interconnect to NVLink 3.0 (600 GB/s per GPU). Variants of Ampere include the A800 (80 GB GDDR6X for China, due to export rules) and computefocused cards like A40, A10, A30, A16 targeting graphics, virtualization, and inference workloads.
- Hopper (2022) NVIDIA's ninth generation (architecture codenamed "Blackwell") debuted with the H100 Tensor Core GPU. H100 is built on the Hopper architecture (≥80 billion transistors on 4 nm) and uses 80 GB HBM3e (3.35 TB/s). It introduced the Transformer Engine for FP8/FP16 mixed-precision (further accelerating large language models ([3] nyidianews.nyidia.com)) and new DPX instructions. Official specs report the H100 SXM module delivers on the order of 67 TFLOPS FP32 and 1.979 TFLOPS FP16 ([9] www.nvidia.com). The H100 also doubled NVLink speed to 900 GB/s per GPU (NVLink 4.0) and supports NVSwitch fabrics for node-scale interconnect. Very recently (2025), NVIDIA has announced the H200 GPU - an incremental Hopper upgrade - now with 141 GB HBM3e and ~4.89 TB/s memory bandwidth, delivering ~241.3 TFLOPS FP16 ([16] www.techradar.com). H200 targets scaling beyond H100 for extreme AI/HPC; multiple H200 cards can also be linked (see below).
- Superchips (Grace) In 2022-2025 NVIDIA combined its GPUs with arm-based CPUs into superchips. The GH200 (Grace + Hopper) and GB300 (Grace + Blackwell) integrate a 96-core Arm CPU with a next-gen GPU on a single package (fabricated on TSMC 5 nm). For example, the GB300 "Grace Blackwell Ultra" features a 72-core Grace CPU and a "Blackwell Ultra" GPU on one die. Its claimed peak Al throughput is ~20 PFLOPS and it supports 784 GB of unified LPDDR5X+HBM3E memory ( www.tomshardware.com). These superchips enable CPU-GPU coherency and are aimed at large AI models and HPC workloads. A key use case is NVIDIA's OVXTM servers and supercomputers where Grace+GPU nodes link via NVLink and InfiniBand.



• Ada Lovelace (2023—) – In 2023, NVIDIA released the *L40*S GPU (codename NVL4 "Ada Lovelace – Data Center") for universal data-center acceleration of Al plus real-time graphics. L40S uses the Ada AD102 GPU (same core as the consumer RTX 4090/Titan Ada) but is configured for data-center use: it has *18,176 CUDA cores, 768 Tensor Cores, 142 RT cores*, and *48 GB GDDR6 with ECC* (<sup>[17]</sup> www.techpowerup.com) (<sup>[3]</sup> nvidianews.nvidia.com). According to NVIDIA, L40S achieves ~5× the single-precision (FP32) throughput of the A100 (<sup>[3]</sup> nvidianews.nvidia.com), and 212 TFLOPS of ray-tracing performance via its RT cores (<sup>[3]</sup> nvidianews.nvidia.com). It sits alongside similar cards (L40 without the S) and provides GPU acceleration for workloads like 3D visualization, virtual workstations, and Al inference, complementing the Hopper GPUs.

This architectural progression led to a **broad spectrum of GPU "solutions"** (accelerator cards and packages) tailored to data-center niches. Table 1 (below) compares the key specs of all current NVIDIA data-center GPUs and Blackwell variants, including Ampere (A100/A800), Blackwell (H100/H200), Ada (L40S), and related products (e.g. workstation GPUs and China-specific chips). Following sections elaborate on each, emphasizing interconnect and real deployments.

| GPU Model                            | Arch.<br>(Node)            | Year | Memory<br>(GB)              | Memory Type   | Mem BW (GB/s)                                | FP32 (TFLOPS)                          | FP16 (Tensor)                              | NVLink/NVSwitch  | MIG Partitions   | TDP<br>(W) |
|--------------------------------------|----------------------------|------|-----------------------------|---------------|--|--|--|--|--|------------|
| NVIDIA A100<br>(SXM)                 | Ampere<br>GA100 (7<br>nm)  | 2020 | 40<br>HBM2e<br>(80*<br>opt) | НВМ2е         | 1555   | ~19.5                                  | 312  | 8-way NVLink (600<br>GB/s per GPU) ( <sup>[2]</sup><br>developer.nvidia.com) | Up to 7 instances ( <sup>[8]</sup> developer.nvidia.com) | 400        |
| NVIDIA A100<br>(PCIe)                | Ampere<br>GA100            | 2020 | 40<br>HBM2e<br>(opt80*)     | HBM2e         | 1555   | ~19.5                                  | 312  | PCIe Gen4 x16 (128<br>GB/s)  | Up to 7  | 250        |
| NVIDIA A800                          | Ampere<br>GA100            | 2021 | 80 GB<br>GDDR6X             | GDDR6X        | ~1215 (eff.)                                 | ~19.5                                  | 312  | Not NVLink-capable<br>(China)  | 1  | 400        |
| NVIDIA H100<br>(SXM)                 | Hopper<br>(Blackwell)      | 2022 | 80 GB<br>HBM3e              | НВМЗе         | 3355   | 67 ( <sup>[9]</sup><br>www.nvidia.com) | 1979 ( <sup>[9]</sup><br>www.nvidia.com)   | 8-way NVLink (900<br>GB/s per GPU) ( <sup>[2]</sup><br>developer.nvidia.com) | Up to 7  | 700        |
| NVIDIA H100<br>(PCIe)                | Hopper<br>(Blackwell)      | 2022 | 80 GB<br>HBM3e              | нвмзе         | 3355   | 60 ( <sup>[9]</sup><br>www.nvidia.com) | 1671 ( <sup>[9]</sup><br>www.nvidia.com)   | 8-way NVLink (900<br>GB/s per GPU) ( <sup>[2]</sup><br>developer.nvidia.com) | Up to 7  | 700        |
| NVIDIA H200<br>(SXM)                 | Hopper<br>(Blackwell)      | 2025 | 141 GB<br>HBM3e             | НВМ3е         | 4890 ( <sup>[16]</sup><br>www.techradar.com) | (n/a)                                  | 241.3 ( <sup>[16]</sup> www.techradar.com) | 8-way NVLink (900<br>GB/s)**   | Up to 7  | ~700       |
| NVIDIA L40S                          | Ada<br>Lovelace<br>(AD102) | 2023 | 48 GB<br>GDDR6<br>w/ECC     | GDDR6         | 864  | ~98.9                                  | 1466 ( <sup>[18]</sup><br>www.nvidia.com)  | PCIe Gen4 x16 (64<br>GB/s)   | -  | 300        |
| NVIDIA RTX<br>Pro 6000D<br>(B40)     | Blackwell-<br>derived      | 2025 | 32 GB<br>GDDR7              | GDDR7         | -  | (n/a)                                  | (n/a)                                      | No NVLink (China)  | -  | ~300       |
| NVIDIA RTX<br>Pro 4000<br>SFF        | Blackwell<br>SFF           | 2025 | -                           | -             | -  | -                                      | -  | -  | -  | -          |
| NVIDIA Tesla<br>V100                 | Volta<br>GV100             | 2017 | 32 GB<br>HBM2               | НВМ2          | 900  | 15.7                                   | 125  | 4-way NVLink (300<br>GB/s)   | _  | 300        |
| NVIDIA Tesla<br>T4                   | Turing<br>TU104            | 2018 | 16 GB<br>GDDR6              | GDDR6         | 300  | 8.1                                    | 65.6                                       | PCIe Gen3 x16 (32<br>GB/s)   | -  | 70         |
| Special:<br>GB200<br>NVL72<br>Design | -                          | 2024 | -                           | -             | -  | -                                      | -  | 72 GPUs @1.8 TB/s<br>each ( <sup>[2]</sup><br>developer.nvidia.com)          | -  | -          |
| Special:<br>GB300<br>"Grace+GPU"     | Grace<br>(ARM) +<br>GPU    | 2025 | 784 GB<br>(total)           | LPDDR5X+HBM3E | 130 TB/s (rack)                              | _                                      | _  | Integrated NVLink<br>(NVL72)   | _  | -          |
| Special: DGX<br>SuperPOD<br>(H100)   | Cluster<br>(NVL72)         | 2024 | -                           | -             | -  | -                                      | -  | 256 GPUs via<br>NVSwitch   | 4 MIG each<br>(x1024vGPU)                                | -          |
| Special:<br>Azure NDv6<br>(GB300)**  | NVL72<br>Cloud<br>Cluster  | 2025 | 37 TB<br>(per<br>rack)      | НВМЗе         | 130 TB/s (per rack)                          | (see text)                             | (see text)                                 | 72 GPUs @1.8 TB/s<br>(per rack)  | -  | -          |
| Special:<br>Nvidia OVX               | DGX-style<br>system        | 2024 | Varied                      | Varied        | -  | -                                      | _  | -  | -  | -          |

Notes: Columns show nominal values. "NVLink domain" lists how many GPUs can intercommunicate (per HGX board or cluster) and link bandwidth. ([2] developer.nvidia.com) ([20] developer.nvidia.com). "MIG" indicates NVIDIA's multi-instance GPU splits (A100/H100 support). "Special" rows denote system-level designs (cluster or superchip) rather than a single card. All values are sourced from NVIDIA references and vendor data ([15] developer.nvidia.com) ([3] nvidianews.nvidia.com) ([2] developer.nvidia.com) ([9] www.nvidia.com).

## **Detailed Characteristics of NVIDIA Datacenter GPUs**

## **Memory and Compute**

NVIDIA's datacenter GPUs distinguish themselves by massive on-board memory and high bandwidth, enabling large models and datasets. For example, the A100 80GB HBM2e version offers 1.555 TB/s memory bandwidth ([15] developer.nvidia.com), while the H100 (80GB HBM3e) doubles that to ~3.355 TB/s ([2] developer.nvidia.com). The newest H200 pushes to ~4.89 TB/s with 141 GB of HBM3e ([16] www.techradar.com). In contrast, graphics-focused cards like the L40S use GDDR6, trading some throughput for lower cost; L40S has 48 GB core with 864 GB/s peak bandwidth ([21] www.techpowerup.com). These differences reflect trade-offs: GDDR6 (on L40S/A40) cuts power and cost but limits inferencing throughput compared to HBM arrays.

Compute performance similarly scales. The A100 delivers ~19.5 TFLOPS of FP32 (64-bit) raw compute (312 TFLOPS at FP16) on the SXM form ([15] developer.nvidia.com). By contrast, H100 SXM reaches 67 TFLOPS FP32 and 1,979 TFLOPS FP16 ([9] www.nvidia.com) - roughly 3.4× and 6× larger. respectively. The H200's increased tensor throughput (241.3 TFLOPS FP16) further amplifies capability (16) www.techradar.com). NVIDIA also reports the L40S ADA card delivers nearly 5× the FP32 throughput of the A100 ([3] nvidianews.nvidia.com), thanks to its greatly expanded shader and tensor core counts (18,176 CUDA cores vs A100's 6,912). Table 1 summarizes these spec-level contrasts. In practice, benchmarks confirm such gaps; for example, ORI Labs notes L40S can match or exceed H100 in many graphics/AI inference tests (given its high shader count) although it lacks NVLink ([18] www.nvidia.com) ([3] nvidianews.nvidia.com).

## Multi-GPU Interconnect: NVLink, NVSwitch, and NVL72

A key component of the datacenter platform is the GPU interconnect. Traditional server architectures have PCIe, but NVIDIA supplements with NVLink (GPU-to-GPU links) and NVSwitch for full mesh fabrics. The Ampere A100/Hopper H100 era uses NVLink 3/4: each H100 GPU can sustain 900 GB/s bidirectional over NVLink ([2] developer.nvidia.com), allowing 8 GPUs on an HGX baseboard to form a single large GPU. NVSwitch chips then stitch multiple NVLink domains together; for instance, in an NVIDIA DGX or SuperPOD, up to 16 or more H100s can communicate as if on one bus via NVSwitch enclosures (enabling all-to-all traffic at~900 GB/s). This is vital for tightly-coupled training tasks. NVIDIA notes that with NVSwitch, up to 256 H100 GPUs can be linked in a single "SuperPOD" (these form the backbone of exascale Al systems) ([22] www.ironsystems.com).

Going beyond NVSwitch, NVIDIA's newest NVL72 architecture shatters that limit. In the NVL72 design (Blackwell GPU version), 72 GPUs occupy one NVLink domain ([2] developer.nvidia.com). Each GPU still enjoys 1.8 TB/s link speed (double H100's) ([2] developer.nvidia.com) – a monumental increase. This design achieves an aggregate AllReduce bandwidth of ~260 TB/s across the rack ([20] developer.nvidia.com). Figure 1 (below) illustrates the difference:

• Figure 1: GPU-to-GPU Network Comparison. The traditional HGX H100 baseboard allows 8 GPUs with 900 GB/s links (NVLink 4.0) ([2] developer.nvidia.com). On the right, the GB200 NVL72 design spans 72 GPUs with 1.8 TB/s links each ([2] developer.nvidia.com) ([20] developer.nvidia.com). (All-Reduce aggregate bandwidth is shown for NVL72.)

| Interconnect          | GPUs per Domain | Per-GPU Bandwidth                               | Aggregate AllReduce BW                           |
|-----------------------|-----------------|---|--|
| NVLink (HGX H100 gen) | 8 GPUs          | 900 GB/s ([2] developer.nvidia.com)             | \~7.2 TB/s                                       |
| NVL72 (GB200 design)  | 72 GPUs         | 1.8 TB/s ( <sup>[2]</sup> developer.nvidia.com) | 260 TB/s ( <sup>[20]</sup> developer.nvidia.com) |

These scaling leaps have profound impact. NVIDIA's analysis shows that moving from 8 to 72 GPUs in an NVLink fabric can accelerate giant AI models by 4-30x ([23] developer.nvidia.com). For example, a GPT-like 1.8 trillion-parameter model ("GPT-MoE-1.8T") could train ~4x faster and serve inference ~30x faster on an NVL72 rack than on 8-GPU systems. The real-world significance is evident in case studies (below) where multi-thousand-GPU clusters rely on NVI 72/NVSwitch fabrics for scale.

## Multi-Instance GPU (MIG) and Virtualization

NVIDIA also designed its GPUs for flexibility under virtualization. Starting with Ampere, GPUs like the A100/H100 can be partitioned by hardware into multiple MIG instances ([8] developer.nvidia.com). For instance, up to 7 separate CUDA instances can run in parallel on one A100, each with independent memory/slice. This boosts utilization in cloud settings. Table 1 notes the maximum MIG splits for each GPU (e.g. "7 instances"). The L40/L40S Ada GPUs do not support MIG (they are pure graphics accelerators), nor do smaller RTX InfiniTerra cards.

#### **Memory Architecture and Efficiency**

Another hallmark of NVIDIA's platform GPUs is advanced memory design. The use of stacked DRAM (HBM2, HBM2e, HBM3e) on NVLink-connected boards effectively merges all GPU memory into one system. For example, an 8-GPU SXM node has 8×40 GB = 320 GB unified memory for an A100 system, or 8×80 GB=640 GB on H100. This coherence is critical for very large models. In the Grace superchips (GB300), the CPU and GPU share a unified address space (with 384-bit LPDDR5X plus HBM to yield hundreds of GBs) ([10] www.tomshardware.com).

Energy and cooling are also constrained by memory and compute. NVIDIA's reference design for NVL72 had to address 120 kW of heat per rack ( [7] developer.nvidia.com). To achieve that, they deployed direct liquid cooling manifolds and specialized blind-mate connectors ([7] developer.nvidia.com). This level of thermal design (7 MW clusters of 72 GPUs) illustrates how pushing datacenter GPU density requires new infrastructure. The OCP-published GB200 reference architecture (in collaboration with Vertiv) even details rack reinforcements and high-current busbars to handle 6,000 lbs of mating force and 1,400 A busbars ([24] pglfmc.com) ([25] developer.nvidia.com). In summary, each GPU generation followed by NVL72 demands commensurate upgrades in power/cooling design, making the GPU a centerpiece of system engineering

# **Case Studies and Deployments**

Azure GB300 NVL72 Supercluster (Microsoft+OpenAI)



In October 2025, Microsoft announced one of the most extreme AI clusters ever built: a **GB300 NVL72 supercomputer** on Azure ([1] www.tomshardware.com). This system stitches together *4,608 NVIDIA Blackwell Ultra GB300 GPUs* across NVL72 racks, each rack containing 72 GPUs and 36 Grace CPUs ([1] www.tomshardware.com). The GPUs are connected with NVLink 5 (1.8 TB/s per GPU) and NVIDIA Quantum-X800 InfiniBand switches. In aggregate, this cluster delivers ~92.1 exaFLOPS of FP4 inference performance ([1] www.tomshardware.com). Even per rack, the performance is staggering: 72 GB300 GPUs plus 36 Grace CPUs yield ~1,440 petaflops and *37 TB* of memory, with 130 TB/s of total memory bandwidth ([1] www.tomshardware.com). According to Microsoft, this cluster specifically accelerates OpenAI training tasks, cutting what used to take months down to weeks. This deployment concretely demonstrates the advantage of NVIDIA's scale: by enabling 72-GPU NVLink domains with GB300, Microsoft can train and serve trillion-parameter models at unprecedented speed ([1] www.tomshardware.com) ([23] developer.nvidia.com).

### **Supercomputer and Cloud Adoptions**

Earlier, supercomputing centers began deploying NVIDIA GPUs at scale. For example, in 2021–22, systems like Indiana University's "Big Red 200" (HPE Cray Shasta) and Germany's Jülich "JUWELS Booster" (Atos) were announced to use NVIDIA A100 GPUs (1121 nvidianews.nvidia.com). These systems leverage 8-GPU nodes with NVLink networks; Perlmutter (NERSC, DOE) similarly combined HPE Shasta with A100 to enable advanced climate and materials simulations (1121 nvidianews.nvidia.com). Even Europe's MareNostrum (BSC) or Sumitomo's River Basin consortium replaced older FPGAs with mixed A100/AMD nodes for AI research. Now in ~2025, leading clouds offer dedicated instances: AWS's P4d and Google's A3 instances used A100s, while newer offerings (AWS P5, Azure NDv6) have H100 or GB300 under the hood. Each provider touts the large FRU counts – e.g. AWS reports thousands of H100s in each P5 rack – effectively giving customers highly scalable clusters on demand.

In the enterprise space, NVIDIA partners have deployed specialized servers. For instance, at NVIDIA's GTC 2023, a new RTX L40S-based server was showcased for virtual workstation and 3D rendering workloads. OEMs (Dell, HPE, Lenovo, Supermicro) now feature L40/L40S in their offerings, optimizing for double-duty Al+graphics tasks ([26] www.nvidia.com) ([3] nvidianews.nvidia.com). Similarly, for on-prem Al deployments, NVIDIA's OVX platform (HVAC-cooled chassis with multiple GPUs) is being adopted by telecom and automotive industries. Pattern-searching startup helps a biotech firm use an H100 cluster for protein folding inference (GPT-type architecture).

At the end-user workstation level, one case is the recently released **Asus ExpertCenter PET900N G3**. This is essentially a workstation built around the *GB300* superchip (<sup>[10]</sup> www.tomshardware.com). It delivers DGX-Station-like performance (20 PFLOPS AI) in a desktop chassis by using the integrated Grace CPU and Blackwell GPU on GB300. This indicates NVIDIA's strategy to proliferate data-center GPUs into smaller form factors via partners, showing the platform's versatility (<sup>[10]</sup> www.tomshardware.com).

#### **Cloud and Enterprise Technologies**

Beyond raw compute, data-center GPUs enable new services. NVIDIA's software stack (CUDA, Merlin for data analytics, Triton Inference Server, etc.) is ubiquitous on these GPUs. OCP contributions like NVIDIA's ConnectX-7 NIC are becoming standards (OCP NIC 3.0) in cloud fabrics. Inference services (e.g. Azure AI, AWS SageMaker endpoints) rely on mixed fleets of GPU accelerators (L40S for vision, H100/H20 for NLP/hyperscale). Even HTTP offloading and video transcoding in data centers is sometimes done on GPUs (commercial servers now include inferred video streams). All these trends underscore how NVIDIA's datacenter GPUs are the cornerstone of modern compute infrastructure.

# **Analysis of NVIDIA GPU Compute and Connectivity**

#### **Performance and Scalability**

We have already highlighted individual GPU FLOPS, but system-level benchmarks also matter. NVIDIA paper data and user reports indicate linear scaling with NVLink: an 8-GPU H100 node scales compute ~8x for large parallel jobs. The NVL72 fabric pushes this further: in principle, 72x scaling is achievable if communication remains hidden by NVLink's high bandwidth (1.8 TB/s) and topology. Analytical models from NVIDIA predict that even as model size growth outpaces single-GPU memory, these fabrics allow model parallelism with modest idle time. For example, the move from 8 to 72 GPUs can yield up to 30x faster inference on massive language models ([23] developer.nvidia.com). In practice, however, such scaling is only attained with optimized software (e.g. fully overlapping communication), which is an active area of research.

Aside from LINPACK-style FLOPS, *Al training metrics* are key. NVIDIA publishes examples: an A100 server (8×A100) can train GPT-3 (175B) ~3× faster than V100 era. Extrapolating, the GB300 NVL72 rack claimed ~92 exaFLOPS in FP4 inference – an unheard-of scale (for reference, the top supercomputers barely started hitting ~20 exaFLOPS total compute in any precision in 2023). This illustrates the trend: as Al model sizes balloon from billions to trillions of parameters, GPU clusters must grow accordingly, a driving factor behind NVL72. It also raises energy concerns: 46× more TFLOPS means 46× more power, which is partly mitigated by GPU efficiency improvements but not fully. Efficiency (performance per watt) has improved generation-to-generation (H100 roughly doubles beyond A100). Still, a full rack at 72×H100 (700W each) plus cooling and CPUs easily exceeds 50 kW, reaching the 120 kW design point (<sup>[7]</sup> developer.nvidia.com).

#### **Market and Ecosystem Context**

NVIDIA's dominance (>98% share (<sup>[4]</sup> www.datacenterdynamics.com)) is not just market size; it reflects a mature ecosystem. Every major cloud offers NVIDIA GPU instances; frameworks from PyTorch/TensorFlow to CUDA SQL analytics and the RAPIDS stack all target these GPUs. The NVIDIA AI Enterprise suite brings Kubernetes APIs and virtualization. Thus, buying NVIDIA is a near-default for AI workloads. Table 1's breadth shows NVIDIA's segmentation: from top-tier H100/H200 for HPC/training, to L40S for inference and visualization, down to T4 for distributed video.

By comparison, AMD's Mi300 or Huawei's Ascend are emerging, but currently limited. For instance, AMD announced its own CDNA-3 Mi300 (6.55 TB/s HBM3) in 2023, but adoption is just beginning. WindowsCentral reports that even with 6x growth in AMD shipments expected ([5] www.reuters.com), NVIDIA retains an order-of-magnitude more units in data centers. One study of Chinese restrictions notes NVIDIA'S China share falling from 95% to ~50% due to export caps ([27] www.reuters.com), signaling that even NVIDIA's edge is challenged by geopolitical forces and native competition (e.g. Huawei).

Nevertheless, the comprehensive feature set (CUDA ecosystem, high bandwidth interconnects, multi-year roadmaps) cements NVIDIA's position.

Experienced HPC practitioners often cite NVIDIA's "software advantage" – a single optimization (CUDA kernels, cuDNN) can run on any future NVIDIA GPU, which is not yet true for competitors. This is reflected in NVIDIA's own statement: the upcoming China-bound GPUs (B30/B40) still rely on NVIDIA's



mature CUDA stack to stay competitive ([5] www.reuters.com). It's clear NVIDIA is preparing downtuned chips (using GDDR7 instead of HBM) to comply with policy while leveraging their ecosystem edge.

# **Future Implications and Directions**

Looking ahead, the NVIDIA data-center platform continues to expand in capability. The **Blackwell architecture** will likely see further iterations (the anticipated H300 beyond H200). NVIDIA's financial filings and press allude to new architectures roughly every 2–3 years. We can expect even larger GPU memory (beyond 141 GB) and more on-chip sparsity functions or new precision modes (e.g. FP8). Integration trends (like Grace networking) suggest GPUs will cooperate even more with DPUs (Data Processing Units) and CPUs. For example, NVIDIA's DOCA software stack is beginning to tie ConnectX NICs and BlueField DPUs into coherent workflows, hinting at future heterogeneous chips or on-silicon converged designs.

The rise of hyperscale GPU clusters also acts as a catalyst for cooling and power innovation. OCP contributions show substantial rack redesign – e.g., 7 MW GB200 NVL72 clusters requiring chilled water at scale ([28] developer.nvidia.com). In data centers, we may see more custom facilities (liquid-cooled pods, immersion cooling) specifically built around GPU densities.

On the application side, NVIDIA's platforms are enabling models of previously impossible scale. Trillion+ parameter LLMs (GPT-4 scale) are now trainable; inference services for them are deployed worldwide. Inference optimizations (FP8, quantization) and specialized runtimes (Triton, TensorRT) will make more efficient use of these GPUs. NVIDIA's roadmap suggests persistent memory (Optane tiers) and other heterogenous memory might be integrated, blurring the lines between GPU and DRAM beyond current unified memory concepts.

Finally, geopolitical factors remain crucial. With China developing homegrown GPUs (e.g. Fenghua 112 GB HBM card [<sup>[29]</sup> www.windowscentral.com)), NVIDIA is likely to produce more segment-specific chips (B-series, A800/H800). Compliance-driven designs (like the RTX Pro 6000D/B40) will proliferate, perhaps creating a more fragmented market. Meanwhile, the US is considering export limits on even wider classes of GPUs – the strategy for NVIDIA is to continue innovating while optimizing designs (e.g. by avoiding banned technologies like HBM) to cover as much demand as possible (<sup>[5]</sup> www.reuters.com) (<sup>[30]</sup> www.tomshardware.com).

Regardless, the technical trajectory is clear: ever-higher computational throughput via denser GPUs and interconnects. This not only advances AI but also drives industries like genomics, physics, climate science, and finance. Farther in the future, NVIDIA is already hinting at successor architectures (they trademarked "NVIDIA Blackwell" and other names beyond Ada), suggesting continuous scaling. The convergence of GPU and CPU (Grace series) may one day yield truly unified compute chips. For now, the NVIDIA datacenter GPU platform – encompassing the spectrum from T4 inference tasks to GB300 superclusters – provides the most potent general-purpose compute power available.

Based on current trends and sources, we conclude that NVIDIA will maintain its leadership into the next decade through both iterative performance gains and strategic ecosystem partnerships. The architecture comparison and data presented here (from vendor documentation and industry analysis ([13] developer.nvidia.com) ([4] www.datacenterdynamics.com)) should serve as a definitive reference on the full range of NVIDIA's data-center GPUs and how they stack up today and evolve tomorrow.

## References

All statements and data above are supported by the sources cited throughout (citations in [brackets]). Key references include NVIDIA's official announcements and technical blogs (15] developer.nvidia.com) (13] developer.nvidia.com) (12] developer.nvidia.com), respected industry news (Tom's Hardware, TechRadar) (11] www.tomshardware.com) (16] www.techradar.com), and market analysis reports (14] www.datacenterdynamics.com) (15] www.reuters.com), as well as vendor specification databases (17] www.techpowerup.com) (19] www.nvidia.com). These sources provide specification tables, performance metrics, deployment case studies, and strategic context for each GPU discussed. Each citation in the text corresponds to a particular line or finding in those sources.

#### **External Sources**

- [1] https://www.tomshardware.com/tech-industry/artificial-intelligence/microsoft-deploys-worlds-first-supercomputer-scale-gb300-nvl72-azure-cluster-4-608-gb300-gpus-linked-together-to-form-a-single-unified-accelerator-capable-of-1-44-pflops-of-inference#:~:Micro...
- [2] https://developer.nvidia.com/blog/nvidia-contributes-nvidia-gb200-nvl72-designs-to-open-compute-project/#:~:Prior...
- [3] https://nvidianews.nvidia.com/news/nvidia-global-data-center-system-manufacturers-to-supercharge-generative-ai-and-industrial-digitalization#:~:rende...
- [4] https://www.datacenterdynamics.com/en/news/nvidia-gpu-shipments-totaled-376m-in-2023-equating-to-a-98-market-share-report/#:~:Nvidi...
- $\label{thm:composition} \begin{tabular}{ll} \hline 151 & https://www.reuters.com/world/china/nvidia-launch-cheaper-blackwell-ai-chip-china-after-us-export-curbs-sources-say-2025-05-24/#:~:Black... \\ \hline \hline 151 & https://www.reuters.com/world/china/nvidia-launch-cheaper-blackwell-ai-chip-china-after-us-export-curbs-sources-say-2025-05-24/#:~:Black... \\ \hline \hline 152 & https://www.reuters.com/world/china/nvidia-launch-cheaper-blackwell-ai-chip-china-after-us-export-curbs-sources-say-2025-05-24/#:~:Black... \\ \hline 152 & https://www.reuters.com/world/chipa/world-chipa-ch$
- [6] https://www.tomshardware.com/pc-components/gpus/nvidia-rtx-pro-6000d-b40-blackwell-gpus-reportedly-set-to-supersede-banned-h20-accelerators-in-china#: ~:The%2...
- $\label{thm:project} \begin{tabular}{ll} \parbox{0.5cm}{$/$} & $/$
- [8] https://developer.nvidia.com/blog/nvidia-ampere-architecture-in-depth/#:~:,perf...
- [9] https://www.nvidia.com/en-us/data-center/h100/#:~:FP32%.
- [10] https://www.tomshardware.com/pc-components/gpus/asus-brings-nvidias-gb300-blackwell-ultra-desktop-superchip-to-workstations-features-up-to-784gb-of-coherent-memory-20-pflops-ai-performance#:~:works...
- [11] https://nvidianews.nvidia.com/news/nvidias-new-ampere-data-center-gpu-in-full-production#:~:The%2...
- $\begin{tabular}{ll} $[12]$ https://nvidianews.nvidia.com/news/nvidias-new-ampere-data-center-gpu-in-full-production\#:$\sim:,ener...$$ and $[12]$ https://nvidianews.nvidia.com/news/nvidias-new-ampere-data-center-gpu-in-full-production#:$\sim:,ener...$$ and $[12]$ https://nvidianews.nvidia.com/news/nvidias-new-ampere-data-center-gpu-in-full-production#:$\sim:,ener...$$ and $[12]$ https://nvidianews.nvidia.com/news/nvidias-new-ampere-data-center-gpu-in-full-production#:$\sim:,ener...$$ and $[12]$ https://nvidianews.nvidias-new-ampere-data-center-gpu-in-full-production#:$\sim:,ener...$$ and $[12]$ https://nvidias-new-ampere-data-center-gpu-in-full-production#:$\sim:,ener...$$ and $[12]$ https://nvidias-gpu-in-full-production#:$\sim:,ener...$$ and $[12]$ https://nvidias-gpu-in-full-production#:$\sim:,ener...$$ and $[12]$ https://nvidias-gpu-in-full-production#:$\sim:,ener...$$
- $\hbox{\tt [14] https://developer.nvidia.com/blog/nvidia-ampere-architecture-in-depth/\#:\sim:spars...}$



- [15] https://developer.nvidia.com/blog/nvidia-ampere-architecture-in-depth/#:~:,arch..
- [16] https://www.techradar.com/pro/huawei-ascend-950-vs-nvidia-h200-vs-amd-mi300-instinct-how-do-they-compare#:~:FP8,A...
- [17] https://www.techpowerup.com/gpu-specs/I40s.c4173#:~:The%2...
- [18] https://www.nvidia.com/en-us/data-center/I40s/#:~:Tenso...
- [19] https://www.tomshardware.com/pc-components/gpus/nvidia-introduces-compact-blackwell-professional-graphics-cards-rtx-pro-4000-sff-and-pro-2000-gpus-laun ched-at-siggraph-2025#:-:2025,...
- [21] https://www.techpowerup.com/gpu-specs/I40s.c4173#:~:L40S%...
- [22] https://www.ironsystems.com/nvidia-solutions/data-center-gpu#:~:strea...
- $[23] \ https://developer.nvidia.com/blog/nvidia-contributes-nvidia-gb200-nvl72-designs-to-open-compute-project/\#: \sim: This \%... + (Annual open contributes open compute-project open compute-project$
- [24] https://pglfmc.com/latest\_news/nvidia-contributes-nvidia-gb200-nvi72-designs-to-open-compute-project-us-pioneer-global-vc-difchq-sfo-india-singapore-riyadh-swiss-our-mind/#:~:NVIDI...
- [26] https://www.nvidia.com/en-us/data-center/I40s/#:~:The%2..
- [27] https://www.reuters.com/world/china/nvidia-launch-cheaper-blackwell-ai-chip-china-after-us-export-curbs-sources-say-2025-05-24/#:~:China...
- [28] https://developer.nvidia.com/blog/nvidia-contributes-nvidia-gb200-nvl72-designs-to-open-compute-project/#:~:match..
- [29] https://www.windowscentral.com/hardware/nvidia/nvidia-designing-nerfed-ai-chip-china-blackwell#:~:2025,...
- [30] https://www.tomshardware.com/pc-components/gpus/nvidia-h20-ai-gpu-inventory-is-limited-but-nvidia-is-making-a-new-b30-model-for-china-to-comply-with-export-restrictions#:~:ln%20...

#### IntuitionLabs - Industry Leadership & Services

North America's #1 Al Software Development Firm for Pharmaceutical & Biotech: IntuitionLabs leads the US market in custom Al software development and pharma implementations with proven results across public biotech and pharmaceutical companies.

Elite Client Portfolio: Trusted by NASDAQ-listed pharmaceutical companies including Scilex Holding Company (SCLX) and leading CROs across North America.

Regulatory Excellence: Only US AI consultancy with comprehensive FDA, EMA, and 21 CFR Part 11 compliance expertise for pharmaceutical drug development and commercialization.

Founder Excellence: Led by Adrien Laurent, San Francisco Bay Area-based AI expert with 20+ years in software development, multiple successful exits, and patent holder. Recognized as one of the top AI experts in the USA.

Custom AI Software Development: Build tailored pharmaceutical AI applications, custom CRMs, chatbots, and ERP systems with advanced analytics and regulatory compliance capabilities.

Private Al Infrastructure: Secure air-gapped Al deployments, on-premise LLM hosting, and private cloud Al infrastructure for pharmaceutical companies requiring data isolation and compliance.

Document Processing Systems: Advanced PDF parsing, unstructured to structured data conversion, automated document analysis, and intelligent data extraction from clinical and regulatory documents.

Custom CRM Development: Build tailored pharmaceutical CRM solutions, Veeva integrations, and custom field force applications with advanced analytics

and reporting capabilities.

Al Chatbot Development: Create intelligent medical information chatbots, GenAl sales assistants, and automated customer service solutions for pharma

companies.

Custom ERP Development: Design and develop pharmaceutical-specific ERP systems, inventory management solutions, and regulatory compliance

Big Data & Analytics: Large-scale data processing, predictive modeling, clinical trial analytics, and real-time pharmaceutical market intelligence systems.

Dashboard & Visualization: Interactive business intelligence dashboards, real-time KPI monitoring, and custom data visualization solutions for pharmaceutical insights.

Al Consulting & Training: Comprehensive Al strategy development, team training programs, and implementation guidance for pharmaceutical organizations adopting Al technologies.

Contact founder Adrien Laurent and team at https://intuitionlabs.ai/contact for a consultation.

#### DISCLAIMER

The information contained in this document is provided for educational and informational purposes only. We make no representations or warranties of any kind, express or implied, about the completeness, accuracy, reliability, suitability, or availability of the information contained herein.

Any reliance you place on such information is strictly at your own risk. In no event will IntuitionLabs.ai or its representatives be liable for any loss or damage including without limitation, indirect or consequential loss or damage, or any loss or damage whatsoever arising from the use of information presented in this document.

This document may contain content generated with the assistance of artificial intelligence technologies. Al-generated content may contain errors, omissions, or inaccuracies. Readers are advised to independently verify any critical information before acting upon it.

All product names, logos, brands, trademarks, and registered trademarks mentioned in this document are the property of their respective owners. All company, product, and service names used in this document are for identification purposes only. Use of these names, logos, trademarks, and brands does not imply endorsement by the respective trademark holders.

IntuitionLabs.ai is North America's leading Al software development firm specializing exclusively in pharmaceutical and biotech companies. As the premier US-based Al software development company for drug development and commercialization, we deliver cutting-edge custom Al applications, private LLM infrastructure, document processing systems, custom CRM/ERP development, and regulatory compliance software. Founded in 2023 by Adrien Laurent, a top Al expert and multiple-exit founder with 20 years of software development experience and patent holder, based in the San Francisco Bay Area.

This document does not constitute professional or legal advice. For specific guidance related to your business needs, please consult with appropriate qualified professionals.

© 2025 IntuitionLabs.ai. All rights reserved.