

# A Review of the Nvidia H100: A New Era of AI Computing

Akhil Khurana

Student Year 13, Hartland International School, Dubai, UAE

---

## ABSTRACT

**Nvidia H100 GPU revolutionizes HPC and AI research by achieving something that was previously thought impossible. Significantly, this advance is architectural, comprising of the addition of Tensor cores and a higher number of CUDA cores thus changing computations in health, and translations among others. This evaluation will explore the details of the H100 built for heavy work tasks like deep learning, scientific simulation, real-time data processing in particular. Comparing these results with the previous chip, it is worth noting two-fold growth in FLOPS of approximately twice and 50 percent increase in memory bandwidth. This justifies its excellence in computation-heavy issues. This is additionally supported by applications that span different areas including medical imaging, driverless vehicles development and computational biology. While economic implication and energy consumption of the H100 GPU may be a concern, they give out excellent investment returns that make them interesting for research-oriented programs and commercial ventures requiring serious computation power. However, more studies are crucial in terms of long term reliability assessments, economic viability aspects, algorithmic optimization for such a contemporary compute giant.**

---

## INTRODUCTION

Artificial intelligence (AI) has been very useful in different industries such as improvements in healthcare and self-driving cars. Artificial intelligence has also contributed greatly in development of chats and robots who help enrich human engagements (Fu et al., n.D ). One of the most significant drivers that promote such improvement is the growing demand for considerable processing power that goes into training artificial intelligence models. Unlike traditional computer processors, which were centrally based in CPU configurations, present-day GPUs are critical when it comes to running massive AI models. As a market leader, Nvidia Corporation constantly pushes the boundaries of what GPUs are capable of doing. At first, they concentrated on gaming by their GeForce series with subsequent development of advanced architectures such as Turing and Ampere for wider coverage beyond general purpose computing. Nvidia's contribution to computational science is profound and influential (Abdelkhalik et al. (2018); NVIDIA Corporation (2020)).

The Nvidia H100 has re-volutionized the field of AI computation. This state-of-the-art GPU not only demonstrates impressive computing capabilities but also aims to re-define the effectiveness and expandability of AI operations. Inspired by Choquette's groundbreaking work (2023), this innovative technology necessitates a thorough evaluation of its real-world applications, performance measurements, and overall impact on various industries it influences.

This paper aims to provide a comprehensive review of the Nvidia H100, such as the architectural innovations, performance benchmarks, and a range of applications from deep learning frameworks. Given the multidisciplinary applications of AI, this review will have a broad scope but will provide particular focus on computationally intensive tasks, where the Nvidia H100's performance could be a necessary factor.

## LITERATURE REVIEW

The first graphics processing unit (GPU) dates back to the 1990's, when they originally aimed at accelerating the graphics. This generation of GPUs was very simple, and its purpose was for lighting and texturing computations. These computers were not designed as the general-purpose type. They could be limited with regard to processing strength, memory capacity, or flexibility. Programmable shaders became available shortly in the early 2000s and even though they allowed for more flexibility, they were still only usable in graphical activities. Researches by researchers show that their efforts to use graphics processors units for other computing tasks gave rise to many challenges, such as memory management, programmability, and power consumption. This transformation of Graphics-Specific CPUs to General-Purpose Parallel Computing Platforms has been largely made possible by

Nvidia. In 2006, the computing unified device architecture (CUDA) platforms introduced by Nvidia greatly enhanced the flexibility of GPUs to handle a large range of computational tasks.

The versions of the company's GPU architectures feature enhancements in the number of CUDA cores with higher memory bandwidth and other optimizations for better parallel processing. The inclusion of technologies such as Tensor Core and Ray Tracing Core has broadened the scope of usability by applications beyond mere graphics rendering. One other example, as mentioned, is the Turing architecture, which has been associated with remarkable advances like real time ray tracing, and AI enhanced visuals, that occurred in 2018 (Fujita, 2022; Skorych, 2022). It broke its previous world record on power consumption. Ampere architecture that is released in the year 2020 is its successor. It took things even further by up raising the number of CUDA cores as well as introducing more advanced Tensor cores tailored specifically for AI and machine learning applications (NVIDIA corporation, 2020). The two architectures are heavily applied in various big data, machine learning, and sciences just demonstrate the trend that GPU is increasingly considered specialist acceleration devices in different hard cases. However, the Nvidia H100, which is essentially the biggest breakthrough among contemporary GPUs that has received immense academic and commercial interest.

There has been a huge number of publications done about it. These include research papers, technical reports, reviews, etc., examining its performance, architecture, and application spectrum. The comprehension drawn from the available literature is that the Nvidia H100 has impressive computing power and displays significant advancements over preceding generations in terms of synthetic benchmarks as well as actual usages. Nonetheless, these reviews identify several issues worth looking into, including the product's long term reliability, cost effectiveness, and capacity to handle complex computations such as in neuroimaging, video creation and other big data related forecasts.

### **Nvidia H100 Architecture**

The Nvidia H100 GPU is based on modern semiconductor fabrication techniques capable to deliver higher power efficiency and computing capacity. Previous versions have more CUDA cores per GPU organized in such a way as to increase parallel processing capability at the same time. Apart from CUDA cores, it has extra Tensor Cores tailor-made to handle memory intensive deep learning tasks. The memory architecture of the GPU comprises the ultrahigh bandwidth ultra-fast GDDR6 memory to ensure fast data transport and storage. Its architecture follows a hierarchy that comprises numerous Streaming Multiprocessors (SMs) containing their own fixed numbers of GPU, CUDA & Tensor cores.

They comprise a high-speed cache hierarchy and memory access pathways with low latency data transfer. The architecture design aims at enhancing performance while minimizing bottlenecks. Therefore, the system is suitable for data-intensive applications such as machine learning, scientific simulations, and image processing. Unlike their predecessors, the Nvidia H100 showcases several new innovations.

Moreover, it makes use of latest Tensor Cores of the upcoming age that have been exclusively developed for speeding up AI processes. The new Tensor Cores are able to accomplish an array of precision matrix multiplication operations much faster compared to their predecessors thus cutting down the time used in training and inference of complex learning models. Also, there has been a marked improvement in architectural design and modern technology to manufacture houses. These optimizations do not only reduce operational expenditures but also make the GPU environmentally friendly and low on global warming emissions.

The H100 also comes with a suite of software tools and SDKs that are created specifically towards simplifying the development and deployment of large AI models. These include improved debugging tools, profiling utilities, and specialized libraries that make it easier for researchers and developers to leverage the full capabilities of the hardware along with the developing AI libraries like TensorFlow and Pytorch.

### **Comparison with Previous Generations**

#### **Benchmark Tests**

To evaluate Nvidia's H100 quantitatively, different benchmark tests have been conducted against previous Nvidia's generations including Turing and Ampere. Detailed examinations were made of various criteria such as calculations using FLOPS, memory bandwidth and task specific latencies.

#### **Comparative Analysis**

Such assessments point out that both global and specialized computing tasks have gone through tremendous changes. For instance, H100 exhibits double increment of FLOPS whereas that for the previous device is only 50%. In addition, when measured inside AI workload's boundary, time spent on training and inference is significantly reduced which makes it compatible for the high powered computational needs.

### **Applications of Nvidia H100**

Nvidia designed the Nvidia H100 architecture with advanced AI models in mind especially due to the next generation of Tensor cores. Some example applications could be natural language processing (NLP), computer vision, and reinforcement learning. It has huge computational throughput which speeds up the training phase for large deep neural networks in the research-development arena. There are many case studies that demonstrate practical applications of H100. H100 Nvidia has enabled creation of self-driven cars that are able to read live sensor info precisely on the streets. In particular, another useful application is found in GAN-based drug development that dramatically speeds up computation times (Bähr et al., 2022; Herrero-Pérez and Martínez Castejón, 2021; Li et al., 2022; Pandey et al., 2022; Skorych and Dosta, 2022).

Likewise, the H100 for various neuroimaging purposes such as voxel-based morphometry; functional MRI analysis; and diffusion tensor imaging. Its high-speed data processing ability has been beneficial in real-time brain mapping that helps doctors arrive at accurate medical diagnosis on time. Neuroimaging has been transformed by the use of convolutional neural networks trained on the H100 which have been applied in many deep learning tasks such as automatic detection of lesions and classification of tumours as well as others.

It is important to note that the H100 has been employed in different fields besides physics including computational biology where it has been useful in projects involving protein-folding simulations, genomic sequence analysis, and molecular dynamics simulations. Its ability to compute helps in running intricate simulations at high speeds thus helping science progress faster in the area.

### **Economic and Environmental Impact**

The cost of running such models at low level is expensive in spite of the fact that running such models at low level is expensive. Businesses and research institutions benefit from cutting down computer consuming activity on such operations because this translates to reduced costs of labor, not forgetting fast product / study findings to market. A substantial ROI can be provided by Nvidia H100 especially firms encountering computational capacity obstacles. AI model training or scientific simulations are accelerated, making organizations monetize their endeavors or disseminate scientific findings at a faster pace (Peddie, 2023b).

The H100 may consume less energy per unit, but compared to previous models, it sure does consume a lot. However, this can prove challenging to data centres or companies seeking carbon neutrality. With so much focus worldwide on reducing carbon dioxide emission, energy conservation is no longer just an economy-related issue, it is an ecology-related one too.

## **DISCUSSION**

As a whole, the various sections of the review provide a detailed account on the architecture, the applications as well as the implication of the Nvidia H100. This machine is equipped with advanced architectural designs as well as new calculation features that may bring about a breakthrough in the sector of AI and supercomputing. Intriguingly, the fascinating applications of the GPUs in AI, deep learning, and medical fields help in highlighting its adaptability and capability. It might be costly but it presents an excellent prospective return on investment with regard to most compute-intensive applications.

It is highly imperative for there to exist standards and laws with regards to these ethical issues concerning GPU's powers since it involves such issues as energy use and algorithmic accountability. This study, however, is limited in scope. Despite trying to cover everything, it's impossible to go into great detail for all topics, like technical specifications and ethics. Furthermore, since the Nvidia H100 is new in the market, its impact on organizations over the long term are yet to be fully studied.

With the launch of the Nvidia H100, the GPUs and computational power houses will have a new benchmark. In addition, future research may focus on different fields like creating new algorithms which are applicable specifically for such powerful devices, refining present-day algorithm's efficiency by utilizing H100 properties.

## **CONCLUSION**

Nvidia's H100 GPU is an enormous leap forward in the field of High Performance Computing (HPC), standing out as being key to developments within Deep Learning and other specialist areas such as medical imaging. This GPU's architecture is very innovative and consists of new-generation tensor cores, which dramatically increase the speed of processing ML tasks. The cores are specially designed to perform tensor operation on a parallel basis, thus enhancing speed of training and inference and efficiency of computations. It also has improved power efficiency, making it a low-cost method of deployment without incurring significant operational expenses related to power needs. This is very crucial, especially in contexts of data centers which often are quite expensive with regards to energy expenses. The diverse applications for the H100 from scientific simulations to operational analytics and

neural networks shows its versatility in bringing about positive change in multiple sectors. In terms of economy, the GPU offers innovative means of developing economic business models and streams of income in so far as it used to carry out calculations which were prohibitively expensive at the time. In an ethical sense, although the H100 offers fantastic prospects for medical breakthroughs, there is a need to adopt a multidisciplinary perspective when involving data privacy and algorithms bias, among other ethical issues. Consequently, the Nvidia H100 is a revolutionary breakthrough in the industry with both challenges and potentials deserving analysis.

## REFERENCES

- [1]. Abdelkhalik, H., Arafa, Y., Santhi, N., Badawy, A.H.A., 2022. Demystifying the Nvidia Ampere Architecture through Microbenchmarking and Instruction-level Analysis, in: 2022 IEEE High Performance Extreme Computing Conference, HPEC 2022. <https://doi.org/10.1109/HPEC55821.2022.9926299>
- [2]. Bähr, P.R., Lang, B., Ueberholz, P., Ady, M., Kersevan, R., 2022. Development of a hardware-accelerated simulation kernel for ultra-high vacuum with Nvidia RTX GPUs. *International Journal of High Performance Computing Applications* 36. <https://doi.org/10.1177/10943420211056654>
- [3]. Choquette, J., 2023. NVIDIA Hopper H100 GPU: Scaling Performance. *IEEE Micro* 43. <https://doi.org/10.1109/MM.2023.3256796>
- [4]. Fu, Y., Bolotin, E., Chatterjee, N., Nellans, D., Keckler, S.W., 2022. GPU Domain Specialization via Composable On-Package Architecture. *ACM Transactions on Architecture and Code Optimization* 19. <https://doi.org/10.1145/3484505>
- [5]. Fujita, K., Yamaguchi, T., Kikuchi, Y., Ichimura, T., Hori, M., Maddegadara, L., 2023. Calculation of cross-correlation function accelerated by TensorFloat-32 Tensor Core operations on NVIDIA's Ampere and Hopper GPUs. *J Comput Sci* 68. <https://doi.org/10.1016/j.jocs.2023.101986>
- [6]. Herrero-Pérez, D., Martínez Castejón, P.J., 2021. Multi-GPU acceleration of large-scale density-based topology optimization. *Advances in Engineering Software* 157–158. <https://doi.org/10.1016/j.advengsoft.2021.103006>
- [7]. Li, B., Patel, T., Samsi, S., Gadepally, V., Tiwari, D., 2022. MISO: Exploiting Multi-Instance GPU Capability on Multi-Tenant GPU Clusters, in: SoCC 2022 - Proceedings of the 13th Symposium on Cloud Computing. <https://doi.org/10.1145/3542929.3563510>
- [8]. NVIDIA Corporation, 2020. NVIDIA A100 Tensor Core GPU Architecture. White Paper.
- [9]. Owens, J.D., Houston, M., Luebke, D., Green, S., Stone, J.E., Phillips, J.C., 2008. GPU computing. *Proceedings of the IEEE* 96. <https://doi.org/10.1109/JPROC.2008.917757>
- [10]. Pandey, M., Fernandez, M., Gentile, F., Isayev, O., Tropsha, A., Stern, A.C., Cherkasov, A., 2022. The transformational role of GPU computing and deep learning in drug discovery. *Nat Mach Intell.* <https://doi.org/10.1038/s42256-022-00463-x>
- [11]. Peddie, J., 2023a. The History of the GPU-Steps to Invention, The History of the GPU-Steps to Invention. <https://doi.org/10.1007/978-3-031-10968-3>
- [12]. Peddie, J., 2023b. The History of the GPU - Eras and Environment, The History of the GPU - Eras and Environment. <https://doi.org/10.1007/978-3-031-13581-1>
- [13]. Skorych, V., Dosta, M., 2022. Parallel CPU–GPU computing technique for discrete element method. *Concurr Comput* 34. <https://doi.org/10.1002/cpe.6839>
- [14]. Van Stigt, R., Swatman, S.N., Varbanescu, A.L., 2022. Isolating GPU Architectural Features Using Parallelism-Aware Microbenchmarks, in: ICPE 2022 - Proceedings of the 2022 ACM/SPEC International Conference on Performance Engineering. <https://doi.org/10.1145/3489525.3511673>