# To What Extent Has Nvidia H100 Heralded A New Age of AI Computing?

**Sarvesh Chandirani**

Year 13 | Jumeirah College, Dubai, UAE

## ABSTRACT

Nvidia H100 GPU has made a revolutionary impact on high power computing (HPC) and AI research. It has achieved previously unthinkable abilities, primarily through architectural improvements by adding more Tensor cores and CUDA cores. This advancement has powerful computational impact in all the fields. This review paper will look deeper into H100 to evaluate its performance in heavy tasks like deep learning, running scientific simulations or processing high quantum of data in real time. A comparison of H100 GPU with the previously available chips reveals a two times increase in FLOPS rates and one and half times increase in memory bandwidth. This substantial increase in performance further increases its significance for some key areas like medical sciences, driverless cars and life sciences. Although the higher costs and energy consumption emerge as potential concerns, these chips promise excellent returns on investment, especially for research and development programs or business ventures that heavily rely on high computational power. However, long term studies will provide further data related to perormance reliability, economic viability and potential impact of optimization of algorithms on the performance of this chip.

## INTRODUCTION

With the recent advancements in the fields of data mining, artificial intelligence and machine learning, our current industrial system has undergone a paradigm shift. This shift is clearly visible in fields like medical diagnostics, self-driving technology, research and development functions and so forth. The interaction between humans and computer is undergoing a massive transformation with the emergence and evolution of intelligent agents like chatbots and robots and these experience are pervasive in the modern world, making them common life experiences (Fu et. Al., n.d.). As artificial intelligence  is powering the ongoing paradigm shift in the modern market and industries, it also requires increasingly larger computational infrastructure to facilitate training of complex AI models that can mimic or emulate certain human cognitive functions.

General purpose based computing tasks are historically handled by CPUs (central processing unit), which have proved to be excellent in carrying out sequential processing of data. However, modern computing needs, especially that of AI based applications, have resulted in increasing demands for parallel computations. These processing demands are particularly enormous when deep learning models are undergoing training process. CPUs have proved to be lacking in processing abilities when handling these demands. This limitation of CPUs resulted in wider application of GPUs (Graphic Processing Units),which have proved to be excellent in handling parallel processing of data. Nvidia Corporation has been a leading innovator in the field of GPUs. Its GeForce GPUs proved to be quite successful among the gaming community. However, the company looked beyond this success to address much broader needs in the fields of research and development, analysis of big data and the development of artificial intelligence.

Advancements in the field of AI is heavily reliant on the evolution of complementary hardware and H100 GPU marks a significant stage in this hardware evolution. This innovative product from Nvidia succeeds the Ampere architecture based A100 GPU chips. Based on Hopper architecture, H100 shows significantly improved processing throughput, memory design (in the context of AI/ML), power efficiency (in the context of its processing performance), and AI framework integrations. With this range of advancements, H100 shows a unique capability to handle the workload demands of AI and real-time data processing and analytics.

Nvidia has redefined industry standards through its focus on constant evolution of GPU architecture (Abdelkhalik et al.2018;  NVIDIA Corporation,2020). H100, at this point of time, stands as a pinnacle of this innovative culture with its unprecedented performance across the spectrum. In the context of AI, these GPUs are expected to cut down the time needed

to train AI models, enhance the accuracy level of inference drawing and provide efficient handling of big data. This paper reviews the capabilities, limitations and potential impact of H100, particularly from AI reliant industries.

### A brief history of GPUs

The basic technology of graphic processing unit (GPU) existed in the 1990s when it was primarily used to render simple graphic content. Limited by rigidity of design and processing capabilities, these GPUs were only good enough for simple lighting and shading functions used in contemporary games and lacked the computational capabilities required for general purposes.

The GPU technology underwent a major transformation in early 2000s with the introduction of shaders that were programmable. While this technological shift offered the computer programmers a greater degree of flexibility, but the use of GPUs was still rather restricted to operations that were graphic intensive. Subsequently, the focus of researches shifted to repurposing GPUs to perform a wider range of computational tasks like running simulations and processing data. However, these research undertakings faced many challenges. These GPUs could be programmed only in limited ways, their memory management was inefficient and consumed high power.

With the introduction of Compute Unified Device Architecture (CUDA) by Nvidia in 2006, the GPU technology underwent a paradigm shifting transformation as the developers were able to use a C like language to create parallel computing programmes. This transformation opened the GPU technology to transcend the traditional role of graphic rendering to a much wider rang eof computational tasks. With the advent of CUDA, GPUs were effectively transformed into high-throughput processors that could handle the demands of scientific researches, commercial functions and AI related tasks.
The subsequent architectures developed and used by Nvidia have made constant improvements through enhanced memory capacity, increased energy efficiency and better ability for parallel processing. With the introduction of Turing architecture that added ray tracing in real time and rendering features with AI enhancements, Nvidia GPUs had become cornderstones of the gaming industry and professional visualization (Fujita, 2022; Skorych, 2022). The Turing architecture was followed by the Ampere architecture in 2020. The CUDA core density was enhanced in this architecture and the Tensor cores were of improved second generation, features that were optimized for the AI industry and deep learning.

This constant evolution in the GPU architecture of Nvidia has made its products virtually indispensable for various fields like climate modelling, researches in the field of genomics or the development of autonomous systems. The advent of Hopper architecture, which underlies the H100 GPUs, offers another game-changing transformation with enhanced capacity of the hardware combined with software support. H100 has been reviewed by various researchers at great lengths, using synthetic benchmarks, and has shown much superior performance as compared to its peer production. This chip has also shown greater compatibility with the range of deep learning frameworks and has cross-sectoral applications like in healthcare, automobile industry and the finance markets.

However, the existing researches also point out some areas of concern such as higher costs, reliability of the hardware in the long term and higher energy consumption. Many studies have pointed out that the depliyent of H100 requires sizeable initial investment and it can make it prohibitively expensive for small startups and research organizations. Other studies that look into the sustainability profile of the technology question its ability to meet the policy targets set for energy efficiency. Irrespective of these concerns, H100 remains a revolutionary product by Nvidia that promises to hold its leading position as a computational hardware. The forthcoming researches and developments in the field of AI infrastructure are expected to revolve around this remarkable piece of technology.

### Looking into the unique architecture of Nvidia H100

The superlative performance of H100 is driven by its Hopper architecture. This GPU is envisioned and built with AI/ML and scientific research in mind. Its key components include 4nm TSMC processors, 80 billion transistors and many other innovative technologies that are first in the industry. All these factors together give a significant boost to its ability to process data in efficient manner.

The unique architecture of H100 uses a higher number of streaming multiprocessors (SMs). Each streaming multiprocessor has many CUDA cores, along with Tensor cores and dedicated caches. These superior built and feature rich streaming processors allow H100 to handle a large number of concurrent threads, thus making it suitable for tasks that are highly data intensive, like training AI models of large scales.

A defining feature of H100 is the tensor cores used in its architecture, which has next generation abilities. These tensor cores show compatibility with multiple formats of data like FP64, T32, BF16 and FP8. This multiformat compatibility allows this GPU superior precision and speed, depending on that task. Its ability to handle FP8 data format is of much significance because it enables delivery of higher throughput but not at the cost of accuracy that is crucial for handling inference tasks.

Another key hardware upgrade for H100 is its HBM3 memory, which comes up to 80GB. Compared to GDDR6 memory, HBM3 offers much higher bandwidth, resulting in significantly higher rate of data transfer and reduction in latency. This higher rate of data transfer with reduced latency plays a crucial role in real-time data processing and running simulations. H100 architecture is further enhanced through the inclusion of NVLink and PCIe Gen 5. NVLinkinterconnects multiple GPUs. These interconnects work at high speed, making them capable of handling massive sets of data or complex models in a seamless manner. While NVLink solves the problem of bottlenecks in data processing, PCIe Gen5 offers maximization of data transfer between GPU and CPU, minimization of idle time and increased efficiency of the overall system.

H100 is further equipped with Multi-instance GPU feature. This is a very interesting feature where a single GPU is transformed into multiple smaller and isolated GPU instances. this capability ensures better allocation of resources, especially in shared environments like data centres or research centres.

H100 has also shown remarkable improvement in efficient use of energy. Although it offers really high level of power, its energy usage even during peak performances is relatively lower. This enhanced efficiency has been achieved through dynamic scaling of voltage and frequency, very sophisticated and advanced cooling mechanism, and optimized allocation of resources through dedicated software.

**A comparative look into benchmark performance**

Benchmark tests are essential for quantitative comparison of H100 with its predecessors and effective estimation of it performance in real-world scenarios. These benchmark tests have been conducted in comparison to preceding architectures, especially A100 based on Ampere architecture and Turing GPUs. These benchmark comparisons look into some key performance indicators- FLOPS , Memory Bandwidth, time taken in AI training and Inference, and Power consumption.

MLPerf benchmark suite was used to conduct these synthetic benchmark tests. These test results showed remarkable improvements in key areas. In assessment of training performance,while training complex AI models like BERT, ResNet and GPT-3, H100 showed 2.5x better performance than A100. With variances based on AI models and the size of dataset, H100 showed improvement of 1.5x to 3x in handling inference based tasks. These improvements have largely resulted from better support for FP8 precision, especially in situations that need rapid model updates or processing data in batches.

H100 also shows 1.5x improvement in memory bandwidth as compared to A100. While A100 had the memory bandwidth of 2 TB/s, H100 shows a remarkable improvement with 3TB/s memor bandwidth. This significant boost in bandwidth allows H100 to handle large datasets with increased efficiency, making it very useful in fields or activities that handle real-time data in large quantities like genomics or trading.

H100 stays ahead of its predecessors during benchmark assessment of its energy efficiency as well. Its per watt performance is higher than the previous generation. This increased energy efficiency leads to lower operational costs and makes it a viable choice for deployment on large scale.

There is significant qualitative improvement in terms of software integration and range of applications. Nvidia AI stack oriented optimization of H100, which includes cuDNN, TensorRT and DeepStream, ensures smoother compatibility with some key frameworks such as TensorFlow, PyTorch and ONNX. Based on the feedback offered by various developers, H100 requires lesser time to setup and model tuning is also easier because it is closely integrated with Nvidia ecosystem.

When tried in real world scenarios that require high-performance computing, H100 In high-performance computing (HPC) environments, the H100 has proven its capability in complex simulations such as fluid dynamics, materials science, and quantum computing. These simulations benefit from both the GPU's raw processing power and its ability to run at higher precision levels when needed.

**Looking into versatility of Nvidia H100 through cross-sectoral applications**

With its market leading innovations in hardware and software, H100 has found a varied range of applications and delivered high quality performance. The architecture of this product facilitates high computational throughput. High computational throughput is a key requirement for the training of any deep learning model, processing of natural language or running complex simulations in the field of research and development.

Training of deep neural networks is a very important field where H100 has wide application. H100 is equipped with advance tensor cores and also has the capability to access data in FP8 model. Both these features combine to reduce the training time required for large language models (LLM) like GPT, BERT or LLaMA. This accelerated training pace means that each iteration cycle is completed at a faster rate. As the iteration cycles are being completed at an accelerated pace, key functions like understanding of natural languages, generative AI and reinforcement of learning can witness faster breakthroughs.

Besides accelerating training time for LLMs, H100 also accelerates the pace at which inference tasks are carried out in AI pipelines. Response of this chip has very low latency, which makes it quite suitable in the field of robotics for making real-time decision, in businesses for providing chatbot based automated customer service, or for providing personalized content on various sites. While performing these tasks, H100 demonstrates efficient handling of simultaneous data streams and provides outcomes at higher speed and reliability. At the same time, resource overhead remains minimal.

H100 has found significant applications in the medical field as well. It has proved to be quite useful in computational biology, discovery of new drugs and neuroimaging. The superior computational ability offered by this chip has accelerated the analysis process of genome sequences. It is quite capable in creating simulations of protein folding and analysis of molecular dynamics. These tasks require high computational power and while the traditional hardware required weeks to complete such tasks, with H100 it can be accomplished in a few hours (Li et al., 2022).

H100 has played a transformative role in the field of neuroimaging as well. It provides analysis of functional MRI with enhanced precision. Other functions like morphometry based on voxel and segmentation of lesion have become more precise. Deep learning models that have been trained usingH100 GPUs show noteworthy enhancement in their ability to detect anomalies in the brain, presence of tumours, and degenerative diseases in their early stage (Pandey et al., 2022; Bähr et al., 2022).

Autonomous vehicle system has also received a major boost with H100. These GPUs are capable of real-time processing of terabytes of data accumulated through various sensors. It allows the vehicles to function efficiently in a dynamic environment by making immediate and accurate decisions. H100 powered AI models are capable of integrating data from diverse sources like camera, LiDAR and navigation system, and quickly process them in real-time to ensure safe and efficient movement.

Besides the automobile sector, H100 has shown applications in robotics and drone based systems as well. H100 enhances the systemic ability to navigate autonomously, detect objects and manipulate in environments that are highly dynamic in nature. The capability of parallel processing at really high speed and creating simulations of path-planning at advanced levels make H100 quite suitable for such systems.

In the financial market there is increasing need for processing market data in real-time and executing highly complex algorithms. H100 has commendable results when applied in the financial market as well. There has been increasing reliance on AI models to detect fraudulent activities, compute credit scoring and perform risk assessment. The speed and accuracy level of these models increases when they are based on the architecture of H100. Nowadays businesses are using these GPUs to analyse consumer sentiments and perform predictive analysis. These applications have shown positive changes in the quality of customer service and the level of operational efficiency.

**A brief look into financial and environmental costs of H100**

While the performance of H100 is market leading, there is a need to evaluate its impact at economic and environmental levels as well.

**financial cost**

Deployment of H100 is capital intensive due its high initial investment. This initial investment includes not only the cost of GPUs but also the supporting infrastructure. This critical infrastructure includes high-performance cooling system to

prevent the GPUs from overheating, better power system and hardware with advanced networking capabilities. However, this high initial investment should not be a prohibitive factor for organizations with high computational power needs. With high computational power the time needed for product development can witness significant reductions, allowing early entry in the market and , thus, significant return on investment (Peddie, 2023b).

While the initial investment for such systems may be prohibitively high for smaller organizations, cloud based access to GPU allowstime-sharing-based use of computing capacity that lowers the cost and makes it possible for smaller organizations to stay competitive against large organizations. This democratizing of computing power creates a more competitive environment that fosters innovation and growth.

### Environmental cost
When we analyse the impact of H100 on the environment, it creates a mixed picture. H100 offer better power efficiency per operations as compared to its predecessors. It would mean that the same task can be accomplished by a smaller number of GPUs. This requirement for lesser number of GPUs can also reduce the power consumption in certain cases. However, when these GPUs are running at max capacity, power consumption goes significantly higher. This poses an environmental concern if these GPUs are deployed in data centres at large scale, because such deployments can result in significant energy consumption.

Since these concerns have been mounting with increasing power costs related to training for AI models and data centres, Nvidia has enhanced the features of H100 to bring power consumption to more acceptable levels. Company has added features like smart scaling of power and optimization of system heating to lower the energy uses. Besides these measures, promoting use of renewable energy to power the H100 based facilities can mitigate the environmental impact. In the future Nvidia can promote circular economy based practices like making the devices more suitable for repair and reuse , promoting recycling of components and , thus, reduce the carbon footprint of these systems in the long run.

### Discussion and conclusion
A cursory review of a selected range of literature clearly establishes H100 as a market leading product that redefines the direction of innovation in the GPU technology. It has many pathbreaking innovations at the level of both hardware and software, like the compatibility with FP8 data format, 2nd generation Tensor cores or the upgrade to HBM3 memory, that establishes it at the forefront of high performance computing. However, it is imperative to look beyond the technical specifics to its transformative impact on the overall system.

The impact of H100 is pervasive, especially in the academic field and commercial sector. Researchers can use the enhanced processing capabilities to try more complex and innovative things in the field of AI like simulation of human cognitive functions or to build AI agents to serve general purposes. Commercial organizations can use the enhanced capabilities to improve the quality of their services and surpass consumer expectations.

However, these benefits are also combined with some serious concerns. The high initial investment required for these systems can prove prohibitive for smaller organizations, leading to artificial disparities in innovation. Although smaller organizations can use cloud based shared access to H100, they may not be competitive against large organizations that can afford large scale deployment of H100s.

There is also an urgent need to explore the ethical aspects of deployment of machine with really high computational capabilities. While high computation power decreases the time needed for training and deploying AI models, it also raises ethical concerns regarding harmful models. While regulatory measures may be put in place, it will also require efforts from Nvidia to introduce technological safeguards against such potential misuses and engage with policymakers to facilitate effective policymaking.

Availability of workforce that is skilled and ready to handle effective and efficient deployment of H100s may also emerge as an area of concern. With all the market leading specifications and new innovations, H100 will require education and training programs to prepare the engineers who canwork with such hardware.

To finally summarize, Nvidia H100 remains a flagship product in the field of GPU technology that can lead to revolutionary changes across a wide-range of fields. However, this transformative impact is predicated upon its deployment in a responsible manner and with a consideration for equitable access.

## REFERENCES

[1]. Y. Choquette et al., "Ten Years of Evolution of NVIDIA GPU Architecture," IEEE Micro, vol. 43, no. 2, pp. 46–55, Mar.–Apr. 2023. [Online]. Available: https://ieeexplore.ieee.org/document/10070122

[2]. NVIDIA, "NVIDIA H100 Tensor Core GPU," NVIDIA Data Center, 2024. [Online]. Available: https://www.nvidia.com/en-us/data-center/h100/

[3]. NVIDIA, "NVIDIA H100 Tensor Core GPU Architecture," NVIDIA Hopper Architecture Resources, 2024. [Online]. Available: https://resources.nvidia.com/en-us-hopper-architecture/nvidia-h100-tensor-c

[4]. Intuition Labs, "The Impact of NVIDIA GPUs in the Pharma Industry," Intuition Labs, 2024. [Online]. Available: https://intuitionlabs.ai/articles/nvidia-gpus-in-pharma-industry

[5]. Uvation, "NVIDIA H100 vs. A100: A Comparative Analysis," Uvation, 2024. [Online]. Available: https://uvation.com/articles/nvidia-h100-vs-a100-a-comparative-analysis

[6]. M. Miller, "NVIDIA H100 GPU Performance Shatters Machine Learning Benchmarks for Model Training," Moor Insights & Strategy, 2022. [Online]. Available: https://moorinsightsstrategy.com/nvidia-h100-gpu-performance-shatters-machine-learning-benchmarks-for-model-training/

[7]. NVIDIA, "NVIDIA Hopper GPU Architecture Technical Overview," White Paper, Version 1.0, 2022. [Online]. Available: https://images.nvidia.com/aem-dam/en-zz/Solutions/data-center/HCC-Whitepaper-v1.0.pdf

[8]. Lambda, "Considerations for Large-Scale NVIDIA H100 Cluster Deployments," Lambda Blog, 2023. [Online]. Available: https://lambda.ai/blog/considerations-for-large-scale-nvidia-h100-cluster-deployments

[9]. eScience Center Denmark, "Denmark's Leap into Advanced AI Research with H100 GPUs," SDU eScience, 2024. [Online]. Available: https://escience.sdu.dk/index.php/news/denmarks-leap-into-advanced-ai-research-with-h100-gpus/

[10]. Princeton University, "Princeton Invests in New 300-GPU Cluster for Academic AI Research," Princeton AI Initiative, 2024. [Online]. Available: https://ai.princeton.edu/news/2024/princeton-invests-new-300-gpu-cluster-academic-ai-research

[11]. NVIDIA, "Confidential Computing on H100 GPUs for Secure and Trustworthy AI," NVIDIA Developer Blog, 2024. [Online]. Available: https://developer.nvidia.com/blog/confidential-computing-on-h100-gpus-for-secure-and-trustworthy-ai/

[12]. Princeton Research Computing, "GPU Computing," Princeton University, 2024. [Online]. Available: https://researchcomputing.princeton.edu/support/knowledge-base/gpu-computing

[13]. Air Street Capital, "Compute Index 2024," Air Street Capital Press, Dec. 2024. [Online]. Available: https://press.airstreet.com/p/compute-index-2024

[14]. IBM, "IBM Brings Enhanced Performance and Efficiency for AI and HPC with NVIDIA Accelerated Computing," IBM Newsroom Blog, 2024. [Online]. Available: https://newsroom.ibm.com/blog-ibm-brings-enhanced-performance-and-efficiency-for-ai-and-hpc-with-nvidia-accelerated-computing

[15]. Harvard Kempner Institute, "Compute Infrastructure," Kempner Institute, 2024. [Online]. Available: https://kempnerinstitute.harvard.edu/compute/

[16]. NVIDIA, "NVIDIA Hopper Architecture In-Depth," NVIDIA Developer Blog, 2022. [Online]. Available: https://developer.nvidia.com/blog/nvidia-hopper-architecture-in-depth/

[17]. NVIDIA, "NVIDIA Hopper Architecture Whitepaper," Advanced Clustering Technologies, 2022. [Online]. Available: https://www.advancedclustering.com/wp-content/uploads/2022/03/gtc22-whitepaper-hopper.pdf

[18]. NVIDIA, "NVIDIA HPC Platform with Hopper and Quantum-2 Achieves Worldwide Adoption," NVIDIA Newsroom, 2022. [Online]. Available: https://nvidianews.nvidia.com/news/nvidia-hpc-platform-hopper-quantum-2-worldwide-adoption

[19]. Paperspace, "NVIDIA's H100: The Powerhouse GPU Revolutionizing Deep Learning," Paperspace Blog, 2023. [Online]. Available: https://blog.paperspace.com/nvidias-h100-the-powerhouse-gpu-revolutionizing-deep-learning/

[20]. Stac Research, "STAC Report: NVDA231030," STAC Research News, 2023. [Online]. Available: https://www.stacresearch.com/news/NVDA231030

[21]. SemiAnalysis, "MI300X vs. H100 vs. H200 Benchmark Part 1: Training," SemiAnalysis, Dec. 2024. [Online]. Available: https://semianalysis.com/2024/12/22/mi300x-vs-h100-vs-h200-benchmark-part-1-training/

[22]. Y. Choquette, M. Gandhi, N. Stam, & R. Krashinsky, "NVIDIA Hopper Architecture In-Depth," IEEE Micro, vol. 43, no. 2, pp. 46–55, 2023. https://doi.org/10.1109/MM.2023.3256796

[23]. Y. Choquette & M. Gandhi, "NVIDIA A100 Tensor Core GPU: Performance and Architecture," Semantic Scholar, 2020. [Online]. Available: https://www.semanticscholar.org/paper/NVIDIA-A100-Tensor-Core-GPU%3A-Performance-and-Choquette-Gandhi/21d0613c3e7fe2cb31f34441c1604edc9882fa45

[24]. Lambda, "Lambda AI," Lambda, 2024. [Online]. Available: https://lambda.ai

[25]. Weights & Biases, "NVIDIA Blackwell GPU Architecture: Unleashing Next-Gen AI Performance," W&B Reports, 2024. [Online]. Available: https://wandb.ai/onlineinference/genai-research/reports/NVIDIA-Blackwell-GPU-architecture-Unleashing-next-gen-AI-performance--VmlldzoxMjgwODI4Mw