

H100 vs. MI300X: A Researcher's Guide to High-Performance GPU Computing

Author: BenchChem Technical Support Team. **Date:** December 2025

Compound of Interest

Compound Name: H100

Cat. No.: B15585327

[Get Quote](#)

In the rapidly evolving landscape of computational research, the choice of graphical processing unit (GPU) can be a critical determinant of a project's success and timeline. For researchers, scientists, and drug development professionals, the NVIDIA **H100** and AMD MI300X represent the pinnacle of GPU acceleration. This guide provides an objective comparison of these two powerful accelerators, supported by experimental data and detailed methodologies, to inform your selection for demanding research applications.

Executive Summary

The NVIDIA **H100**, built on the Hopper architecture, and the AMD MI300X, powered by the CDNA 3 architecture, are both formidable tools for scientific discovery. The **H100**'s strength lies in its mature and extensive CUDA software ecosystem, which offers a seamless experience for a wide array of scientific applications. In contrast, the MI300X boasts a significant advantage in memory capacity and bandwidth, making it particularly well-suited for handling massive datasets and large models. While the **H100** often shows superior performance in highly optimized, production-level AI workloads, the MI300X demonstrates competitive and sometimes superior performance in memory-bound tasks and at larger batch sizes. The choice between them will ultimately depend on the specific requirements of your research, including the nature of your computational tasks, your reliance on existing software, and the scale of your data.

Data Presentation: A Comparative Overview

The following tables summarize the key specifications and performance metrics of the NVIDIA H100 and AMD MI300X.

Table 1: Architectural and Memory Specifications

Feature	NVIDIA H100 (SXM5)	AMD Instinct MI300X
GPU Architecture	Hopper	CDNA 3
Manufacturing Process	TSMC 4N	TSMC 5nm & 6nm
Transistors	80 Billion	153 Billion[1]
GPU Memory	80 GB HBM3	192 GB HBM3[1][2][3]
Memory Bandwidth	3.35 TB/s[1][2]	5.3 TB/s[1][2]
L2 Cache	50 MB	256 MB Infinity Cache[4]
Interconnect	4th Gen NVLink (900 GB/s)	4th Gen Infinity Fabric
Max Power Consumption	700W	750W

Table 2: Peak Theoretical Performance

Precision	NVIDIA H100 (SXM5)	AMD Instinct MI300X
FP64 (Double Precision)	67 TFLOPS	81.7 TFLOPS
FP32 (Single Precision)	67 TFLOPS	163.4 TFLOPS
TF32 Tensor Core	989 TFLOPS	1,305 TFLOPS
FP16/BF16 Tensor Core	1,979 TFLOPS	2,610 TFLOPS
FP8 Tensor Core	3,958 TFLOPS	5,229 TFLOPS
INT8 Tensor Core	3,958 TOPS	5,229 TOPS
*With Sparsity		

Performance Benchmarks and Experimental Protocols

Direct, peer-reviewed performance comparisons in a wide range of scientific applications are still emerging. Much of the available data focuses on Large Language Model (LLM) inference and training, which can serve as a proxy for other computationally intensive tasks.

Large Language Model (LLM) Inference

LLM inference is a memory-bandwidth-intensive task, and the MI300X's superior memory specifications often translate to a performance advantage, particularly with large models and batch sizes.

Experimental Protocol: LLM Inference Throughput

- Objective: To measure the inference throughput (tokens per second) of the **H100** and MI300X on a large language model.
- Model: Mixtral 8x7B.[5]
- Hardware:
 - NVIDIA: 8x **H100** SXM5 GPUs with NVLink.[6]
 - AMD: 8x MI300X accelerators.[6]
- Software:
 - NVIDIA: CUDA 12.2, vLLM v4.3.[6]
 - AMD: ROCm 6.1.2, MK1's inference engine (Flywheel) v0.9.2, and AMD's ROCm optimized fork of vLLM v0.4.0.[6]
- Methodology:
 - The Mixtral 8x7B model was loaded onto the respective GPU platforms. Due to its size, tensor parallelism was set to 2 for the **H100**, while the MI300X could accommodate the entire model on a single GPU.[5]

- Offline inference performance was measured across a range of batch sizes from 1 to 1024.[6]
- Throughput was recorded as the number of tokens generated per second.

Results Summary: In these tests, the MI300X demonstrated a significant performance uplift, ranging from 22% to nearly 3x that of the **H100** as batch sizes increased.[6] At smaller batch sizes, the **H100** has been shown to have a throughput advantage in some tests.[5]

High-Performance Computing (HPC)

For traditional HPC workloads, both GPUs offer substantial double-precision (FP64) performance. The MI300X has a higher theoretical peak FP64 performance.

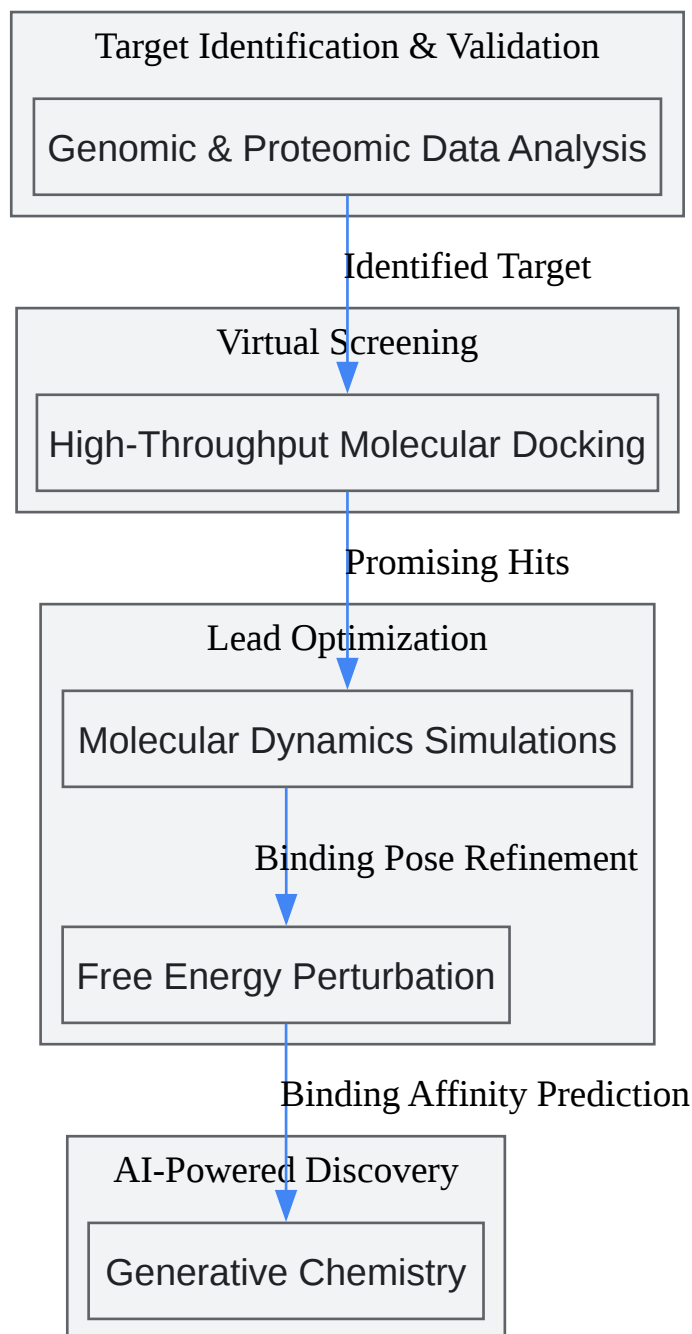
Experimental Protocol: Fast Fourier Transforms (FFT)

- Objective: To compare the memory bandwidth utilization of the **H100** and MI300X in a memory-bound HPC task.
- Benchmark: 1D batched power of 2 complex-to-complex FFTs in single and double precision.
- Software: VkFFT, cuFFT, and rocFFT libraries.[7]
- Methodology:
 - The benchmark was run on both single and double precision.
 - Estimated bandwidth was calculated as $(2 * \text{System size [GB]}) / \text{execution time [s]}$. [7][8]
This metric reflects the efficiency of data transfer between VRAM and the compute units. [7]

Results Summary: In single precision, both GPUs achieved similar results, approaching 3TB/s. [7] In double precision, the MI300X achieved a higher base bandwidth.[7] Neither GPU reached its theoretical peak memory bandwidth, which is common in real-world applications.[7][8]

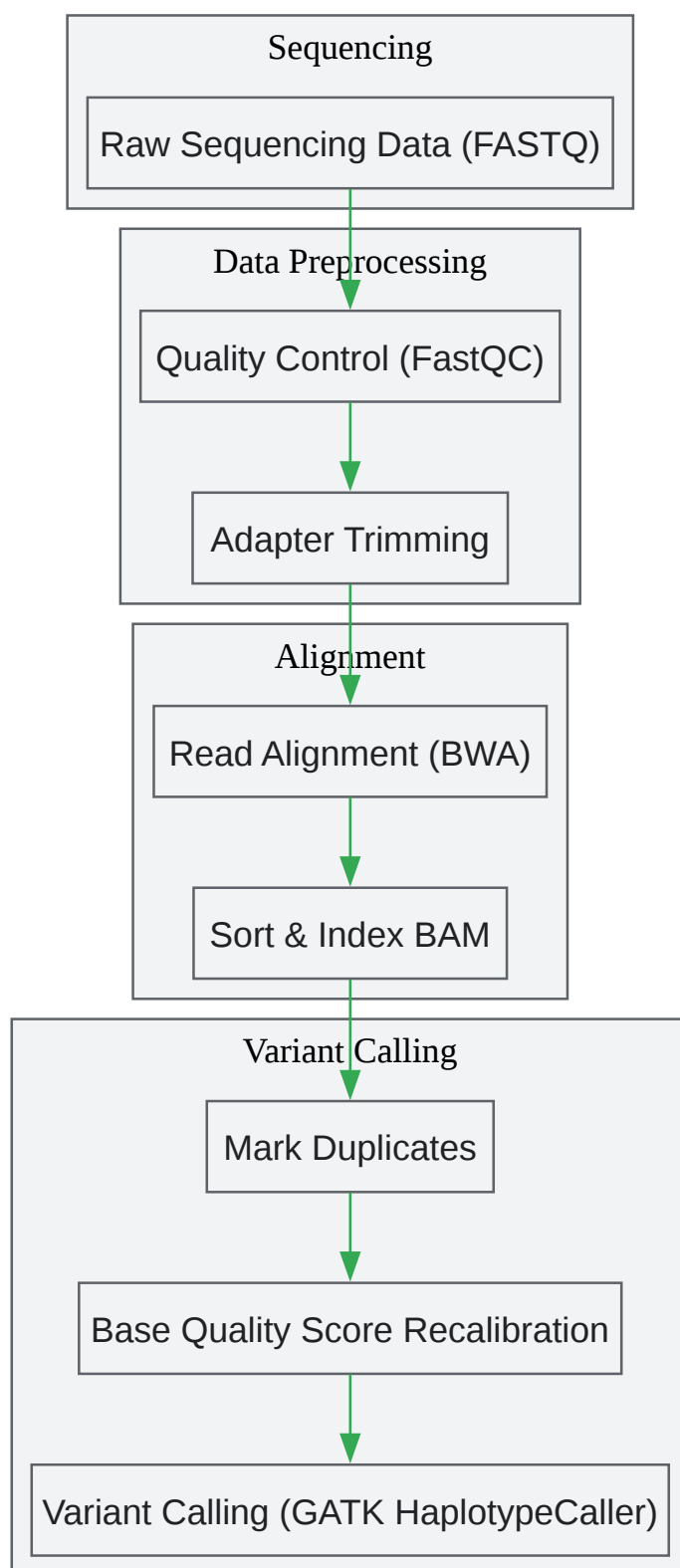
Visualizing Research Workflows

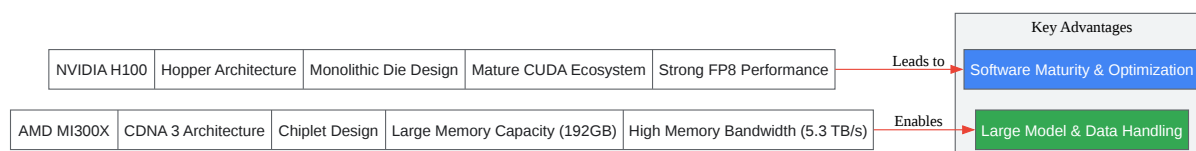
The following diagrams, generated using Graphviz, illustrate common experimental workflows in drug discovery and genomics where these GPUs can be applied.



[Click to download full resolution via product page](#)

Caption: A GPU-accelerated drug discovery workflow.





[Click to download full resolution via product page](#)

Need Custom Synthesis?

BenchChem offers custom synthesis for rare earth carbides and specific isotopic labeling.

Email: info@benchchem.com or [Request Quote Online](#).

References

- 1. AMD Radeon Instinct MI300X Specs | TechPowerUp GPU Database [techpowerup.com]
- 2. Accelerating Drug Discovery: GPU-Enhanced Computational Biology Methods for Molecular Docking Simulations and Virtual Screening[v1] | Preprints.org [preprints.org]
- 3. Accelerate Drug Discovery with GPU-Powered HPC [siliconmechanics.com]
- 4. amd.com [amd.com]
- 5. AMD MI300X vs. Nvidia H100 SXM: Performance Comparison on Mixtral 8x7B Inference | Runpod Blog [runpod.io]
- 6. wccftech.com [wccftech.com]
- 7. GitHub - AI-Hypercomputer/gpu-recipes: Recipes for reproducing training and serving benchmarks for large machine learning models using GPUs on Google Cloud. [github.com]
- 8. reddit.com [reddit.com]
- To cite this document: BenchChem. [H100 vs. MI300X: A Researcher's Guide to High-Performance GPU Computing]. BenchChem, [2025]. [Online PDF]. Available at: [https://www.benchchem.com/product/b15585327#comparing-h100-and-mi300x-for-research-applications]

Disclaimer & Data Validity:

The information provided in this document is for Research Use Only (RUO) and is strictly not intended for diagnostic or therapeutic procedures. While BenchChem strives to provide accurate protocols, we make no warranties, express or implied, regarding the fitness of this product for every specific experimental setup.

Technical Support: The protocols provided are for reference purposes. Unsure if this reagent suits your experiment? [[Contact our Ph.D. Support Team for a compatibility check](#)]

Need Industrial/Bulk Grade? [Request Custom Synthesis Quote](#)

BenchChem

Our mission is to be the trusted global source of essential and advanced chemicals, empowering scientists and researchers to drive progress in science and industry.

Contact

Address: 3281 E Guasti Rd
Ontario, CA 91761, United States
Phone: (601) 213-4426
Email: info@benchchem.com