

# Technical Support Center: Optimizing CUDA Code for the NVIDIA H100

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

## Compound of Interest

Compound Name: H100

Cat. No.: B15585327

[Get Quote](#)

Welcome to the technical support center for optimizing CUDA applications on the NVIDIA **H100** Tensor Core GPU. This guide provides troubleshooting advice, frequently asked questions (FAQs), and best practices tailored for researchers, scientists, and drug development professionals who are leveraging the **H100** for their computational experiments.

## Frequently Asked Questions (FAQs)

### General Performance & Optimization

Q1: I'm migrating my CUDA application from a previous generation GPU (e.g., A100) to the **H100**. What are the first steps I should take to optimize my code?

Applications that followed best practices for previous NVIDIA architectures like Ampere should generally see a performance increase on the **H100** without code changes.<sup>[1][2][3]</sup> However, to unlock the full potential of the Hopper architecture, consider these initial steps:

- **Update Software Stack:** Ensure you are using the latest NVIDIA drivers, a compatible CUDA Toolkit (version 11.8 or later is required for **H100**), and updated versions of cuDNN and other CUDA-X libraries.<sup>[4][5]</sup>
- **Recompile with Hopper Architecture Support:** When compiling your CUDA kernels, target the Hopper architecture by including the appropriate -gencode flags for compute capability 9.0 (e.g., -gencode arch=compute\_90,code=sm\_90). This ensures your application uses the native **H100** instruction set.<sup>[6]</sup>

- **Profile Your Application:** Use NVIDIA's Nsight suite of tools to establish a performance baseline. Nsight Systems will provide a high-level overview to identify bottlenecks, and Nsight Compute will allow for in-depth kernel analysis.[\[7\]](#)[\[8\]](#)[\[9\]](#)
- **Leverage Tensor Cores with FP8 and Transformer Engine:** For AI and deep learning workloads, especially those involving transformer models, the **H100** introduces the FP8 data type and the Transformer Engine.[\[10\]](#) The Transformer Engine intelligently manages and dynamically selects between FP8 and 16-bit calculations to maximize speed while preserving accuracy.[\[10\]](#)

Q2: My application is not performing as expected on the **H100**. What are the common performance bottlenecks?

Performance issues can stem from various factors. A systematic approach to identifying the bottleneck is crucial.[\[4\]](#)

- **GPU Underutilization:** If your GPU utilization is low (monitor with `nvidia-smi`), the bottleneck is likely outside the GPU.[\[4\]](#)
  - **Data Loading:** Slow data pipelines can starve the GPU.[\[4\]](#) Consider using faster storage like NVMe SSDs, enabling data prefetching, and optimizing data preprocessing steps.[\[4\]](#) GPU Direct Storage (GDS) can also help by bypassing the CPU.[\[11\]](#)
  - **CPU-Bound:** The CPU may not be able to submit work to the GPU fast enough. Look for high CPU utilization. Using CUDA Graphs can reduce kernel launch overhead.[\[12\]](#)[\[13\]](#)
- **Memory Bandwidth Limitations:** The application might be limited by the speed at which data can be moved.
  - **Inefficient Global Memory Access:** Ensure global memory accesses are "coalesced," meaning threads in a warp access contiguous memory locations.[\[1\]](#)[\[14\]](#)
  - **Data Transfers:** Minimize data transfers between the host (CPU) and the device (GPU).[\[1\]](#)
- **Kernel Inefficiencies:** The CUDA kernel itself might be suboptimal. Profiling with Nsight Compute can reveal issues like instruction latency, shared memory bank conflicts, or inefficient thread block configurations.[\[8\]](#)

- Thermal Throttling: The **H100** can generate significant heat.[4] Monitor GPU temperatures using nvidia-smi. Ensure proper server cooling and airflow to prevent the GPU from reducing its clock speeds.[4][7]

Q3: How can I effectively utilize the **H100**'s 4th Generation Tensor Cores?

The **H100**'s Tensor Cores are specialized hardware units that dramatically accelerate matrix operations, which are fundamental to AI and many HPC workloads.[14][15]

- Use Mixed Precision: Leverage lower precision data types like FP16, BF16, and the new FP8 format.[10][13] Tensor Cores operate on these data types much faster than on FP32 or FP64.[16]
- Align Matrix Dimensions: Ensure that the dimensions of your matrices are multiples of 16 for optimal Tensor Core performance.[17]
- Utilize CUDA Libraries: Libraries like cuBLAS and cuDNN are optimized to use Tensor Cores automatically for supported operations.
- Program with WMMA/MMA: For custom kernels, use the Warp-Level Matrix-Multiply-Accumulate (WMMA) API or the mma.sync instruction in PTX to directly program the Tensor Cores.[17]

## Memory Optimization

Q4: How should I optimize for the **H100**'s memory hierarchy?

The **H100** features a sophisticated memory hierarchy consisting of HBM3 memory, a large L2 cache, and per-SM shared memory.[12][17][18] Effective utilization is key to performance.[18]

- Maximize Data Locality:
  - Shared Memory: Use the fast, on-chip shared memory to reduce trips to global memory. Tiling or blocking is a common technique where a larger problem is broken into smaller chunks that fit into shared memory.[14][17] The **H100** increases the shared memory capacity per SM to 228 KB.[2]

- L2 Cache: The **H100** has a 50MB L2 cache.[\[12\]](#)[\[16\]](#) It helps to cache frequently accessed data, reducing latency.[\[12\]](#) While you don't control it directly, access patterns that reuse data will benefit.[\[17\]](#)
- Efficient Global Memory (HBM3) Access: The **H100** uses HBM3 memory with up to 3 TB/s of bandwidth.[\[12\]](#)[\[13\]](#) To leverage this, ensure your memory accesses are coalesced.
- Asynchronous Data Transfers: Use the Tensor Memory Accelerator (TMA) to perform asynchronous data transfers between global and shared memory, which can hide data movement latency behind computation.[\[12\]](#)[\[19\]](#)

## Quantitative Data Summary

For a clear comparison, the table below summarizes the key specifications of the NVIDIA **H100** GPU against its predecessor, the A100.

Feature	NVIDIA A100 (SXM4)	NVIDIA H100 (SXM5)	Improvement Factor
Architecture	Ampere	Hopper	-
CUDA Cores	10,752	16,896[13]	~1.6x
Tensor Cores	3rd Generation	4th Generation[20]	-
Max FP8 Tensor Core Performance	N/A	4,000 TFLOPS	-
Max FP16 Tensor Core Performance	624 TFLOPS	2,000 TFLOPS	~3.2x
Max FP64 Performance	19.5 TFLOPS	67 TFLOPS	~3.4x
GPU Memory	Up to 80GB HBM2e	80GB HBM3[13]	-
Memory Bandwidth	2,039 GB/s	3,350 GB/s	~1.6x
L2 Cache	40 MB	50 MB[2][16]	1.25x
NVLink	3rd Generation (600 GB/s)	4th Generation (900 GB/s)[10][13]	1.5x
Max Thermal Design Power (TDP)	400W	700W[10]	1.75x

## Troubleshooting Common Issues

Q5: My GPU utilization is consistently low. How do I troubleshoot this?

Low GPU utilization indicates that the GPU is often waiting for data or instructions.

- Use a Profiler: Run your application with NVIDIA Nsight Systems.[9] The timeline view will clearly show gaps in GPU execution and highlight if the CPU is the bottleneck (CPU-bound) or if data transfers are taking too long (I/O-bound).
- Check for CPU Bottlenecks: If the profiler shows the CPU is at 100% while the GPU is idle, your application is CPU-bound.

- Solution: Offload more parallel parts of the workload to the GPU. Use CUDA Graphs to minimize the overhead of launching many small kernels.[\[13\]](#)
- Inspect Data Pipeline: If the GPU is waiting on I/O, the profiler timeline will show large gaps filled with data transfer operations.
  - Solution: Optimize your data loading process. Use faster storage, implement data prefetching, or use libraries designed for efficient data loading.[\[4\]](#)
- Small Workloads: If you are running many small, independent tasks, the kernel launch latency can dominate the execution time.
  - Solution: Batch smaller tasks into larger kernel launches to improve efficiency.

Q6: My application performance is degrading over time, and the GPU temperature is high. What's happening?

This is a classic sign of thermal throttling. The GPU automatically reduces its clock speed to prevent overheating when it exceeds a certain temperature limit.[\[4\]](#)[\[7\]](#)

- Monitor Temperature: Use the `nvidia-smi` command-line tool to check the GPU's current temperature and power draw.[\[7\]](#)
- Improve Cooling: Ensure the server chassis has adequate airflow and that the data center's cooling systems are functioning correctly.[\[4\]](#)[\[7\]](#) Dust buildup on heatsinks can also be a cause, so ensure hardware is clean.
- Check Power Delivery: Insufficient or unstable power can also lead to performance issues. Verify that the power supply meets the **H100**'s requirements (up to 700W for the SXM5 variant).[\[5\]](#)[\[10\]](#)

## Experimental Protocols

### Protocol: Performance Bottleneck Analysis

This protocol outlines a systematic approach to identifying and addressing performance bottlenecks in a CUDA application on the **H100**.

Objective: To identify the primary performance limiter (e.g., memory bandwidth, compute, latency) in a CUDA application and measure the impact of a targeted optimization.

Methodology:

- Establish a Baseline:
  - Compile the unoptimized application with Hopper architecture support (-gencode arch=compute\_90,code=sm\_90).
  - Run the application on a dedicated **H100** GPU.
  - Use nvidia-smi to log GPU utilization, memory usage, temperature, and power draw during the run.
  - Profile the application using NVIDIA Nsight Systems to get a high-level trace of CPU-GPU interaction and identify major time sinks.[\[9\]](#)[\[21\]](#) Record the total runtime.
- High-Level Analysis (Nsight Systems):
  - Examine the Nsight Systems timeline.
  - Identify: Long periods of data transfer between host and device, gaps between kernel executions (indicating CPU-side bottlenecks), or specific kernels that consume the majority of the GPU time.
- Kernel-Level Analysis (Nsight Compute):
  - Based on the Nsight Systems report, select the most time-consuming kernel for in-depth analysis.
  - Re-run the application profiled with NVIDIA Nsight Compute.[\[8\]](#)[\[9\]](#)
  - Analyze Key Metrics:
    - Memory Throughput: Compare the achieved global memory throughput to the theoretical maximum of the **H100**. Low throughput suggests inefficient memory access patterns.

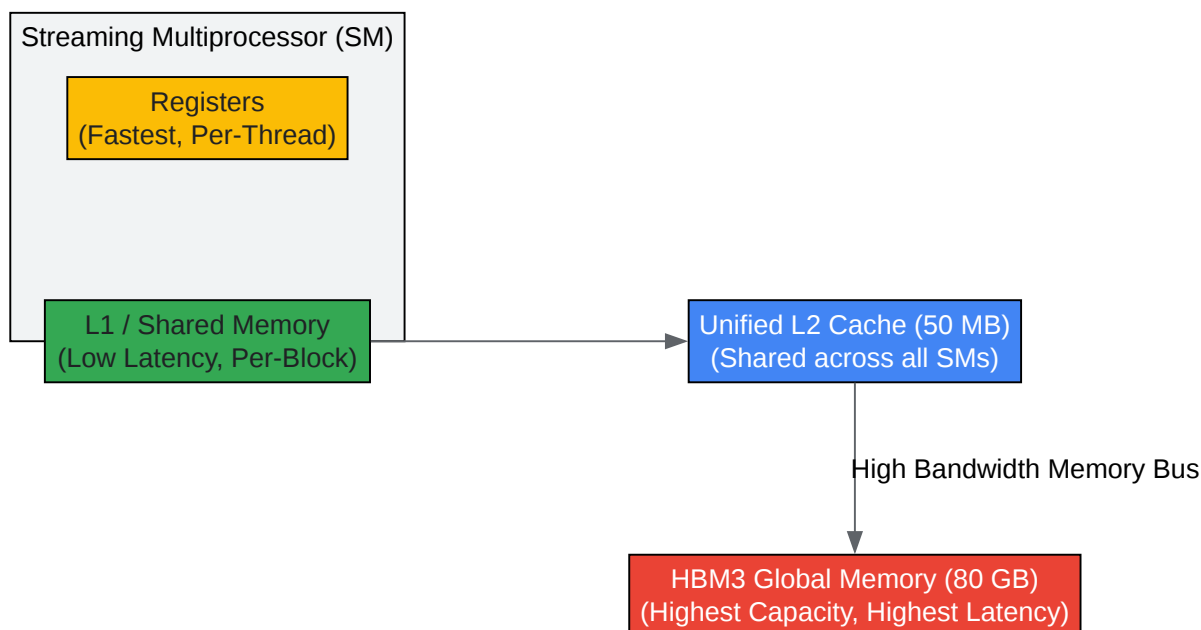
- SM Occupancy: Check the ratio of active warps per SM to the maximum supported.[18]  
Low occupancy can indicate that not enough parallelism is being exposed.
- Instruction Mix: Look at the breakdown of instructions. A high number of memory instructions relative to compute instructions suggests a memory-bound kernel.
- Hypothesize and Implement Optimization:
  - Based on the analysis, form a hypothesis. For example, "The kernel is memory-bound due to uncoalesced global memory accesses."
  - Implement a single, targeted code change to address this. For the example, this would involve refactoring the memory access pattern.
- Measure and Verify:
  - Re-run the profiling steps (1-3) with the optimized code.
  - Compare the new runtime and profiling metrics to the baseline.
  - Verify: Did the runtime decrease? Did the targeted metric (e.g., memory throughput) improve? Did the optimization introduce any new bottlenecks?
- Iterate: Continue this process, addressing one bottleneck at a time, until performance goals are met or no further significant improvements are observed.

## Visualizations

### H100 GPU Memory Hierarchy

The following diagram illustrates the different levels of memory available to CUDA threads running on an **H100** Streaming Multiprocessor (SM). Accessing memory higher up in the diagram (closer to the SM) is significantly faster.



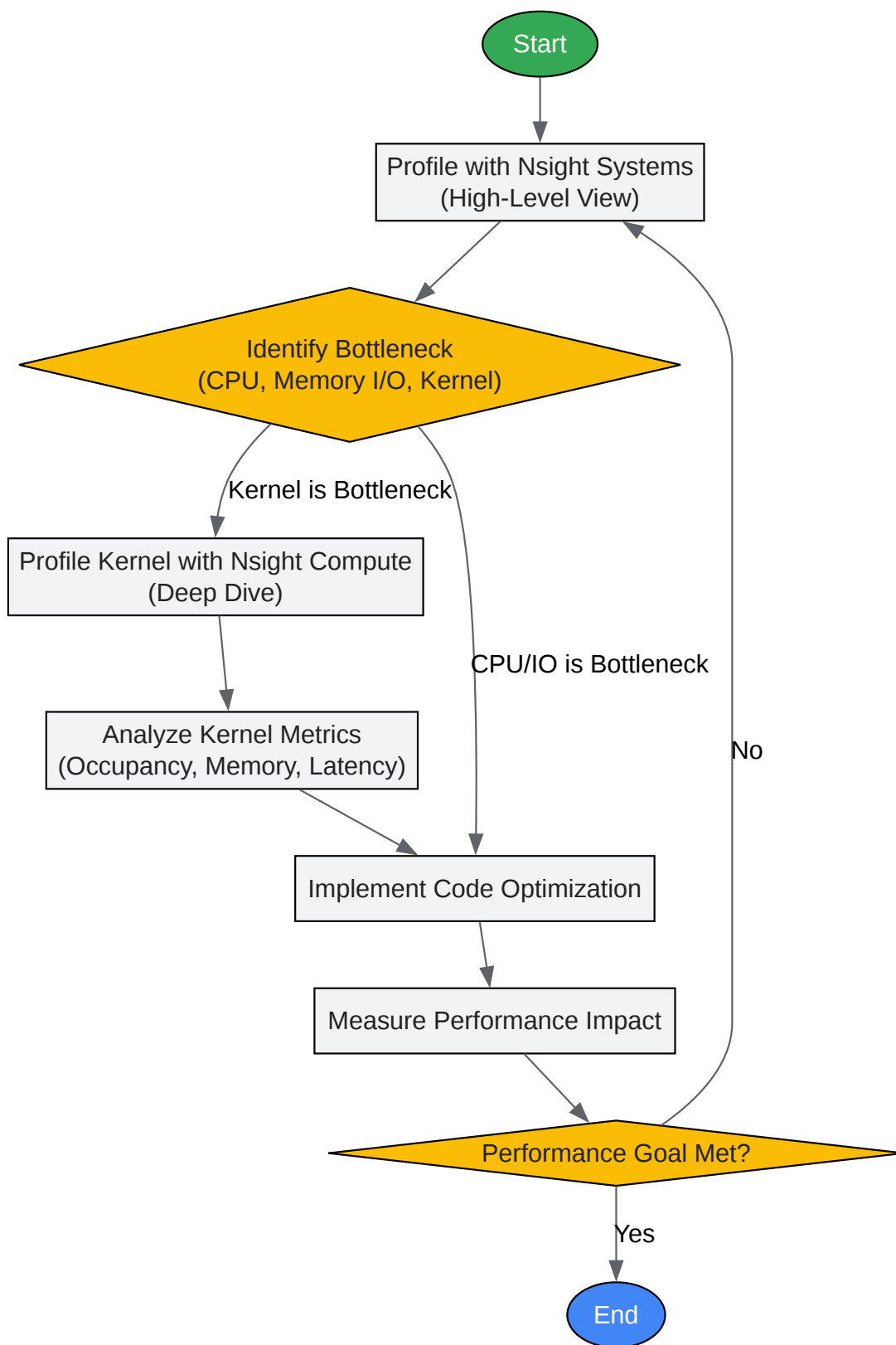


[Click to download full resolution via product page](#)

Caption: Logical diagram of the NVIDIA **H100** memory hierarchy.

## General CUDA Performance Optimization Workflow

Optimizing a CUDA application is an iterative process of profiling, analyzing, and refining. The workflow below outlines the recommended steps.



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

Caption: An iterative workflow for CUDA performance optimization.

### Need Custom Synthesis?

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

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

## References

- 1. docs.nvidia.com [docs.nvidia.com]
- 2. docs.nvidia.com [docs.nvidia.com]
- 3. docs.nvidia.com [docs.nvidia.com]
- 4. How do I troubleshoot common performance issues with the NVIDIA H100 GPU during large language model training? - Massed Compute [massedcompute.com]
- 5. How do I troubleshoot common issues with NVIDIA H100 GPU performance in a distributed training setup? - Massed Compute [massedcompute.com]
- 6. 1. Hopper Architecture Compatibility & BenchChem Hopper Compatibility Guide 13.1 documentation [docs.nvidia.com]
- 7. How do I troubleshoot common issues with NVIDIA H100 GPUs in a data center? - Massed Compute [massedcompute.com]
- 8. Can NVIDIA Nsight Compute be used to optimize GPU-accelerated applications on NVIDIA data center GPUs such as the A100 or H100? - Massed Compute [massedcompute.com]
- 9. 17.3. GPU Profiling — Kempner Institute Computing Handbook [handbook.eng.kempnerinstitute.harvard.edu]
- 10. lambda.ai [lambda.ai]
- 11. How do I troubleshoot common issues with H100 GPU performance in a large-scale HPC cluster? - Massed Compute [massedcompute.com]
- 12. What are the optimal settings for the H100 GPU's memory hierarchy in large language model training? - Massed Compute [massedcompute.com]
- 13. cyfuture.cloud [cyfuture.cloud]
- 14. How do I optimize my CUDA code for NVIDIA A100 and H100 GPUs? - Massed Compute [massedcompute.com]
- 15. medium.com [medium.com]

- 16. medium.com [medium.com]
- 17. Can you provide an example of how to optimize the memory hierarchy of the H100 PCIe GPU for a specific dense matrix multiplication workload? - Massed Compute [massedcompute.com]
- 18. How does the H100 GPU's memory hierarchy affect CUDA occupancy? - Massed Compute [massedcompute.com]
- 19. youtube.com [youtube.com]
- 20. What is an NVIDIA H100? | DigitalOcean [digitalocean.com]
- 21. How do I use NVIDIA's Nsight Systems to profile memory usage on A100 and H100 GPUs? - Massed Compute [massedcompute.com]
- To cite this document: BenchChem. [Technical Support Center: Optimizing CUDA Code for the NVIDIA H100]. BenchChem, [2025]. [Online PDF]. Available at: [https://www.benchchem.com/product/b15585327#optimizing-cuda-code-for-the-nvidia-h100]

---

#### 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](mailto:info@benchchem.com)