

H100 GPU Technical Support Center: Troubleshooting Memory Allocation Errors

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 troubleshooting memory allocation errors on NVIDIA **H100** GPUs. This resource is designed for researchers, scientists, and drug development professionals to diagnose and resolve common memory-related issues encountered during their experiments.

Troubleshooting Guides

This section provides detailed walkthroughs for identifying and resolving specific memory allocation errors.

Issue: "CUDA out of memory" Error During Model Training

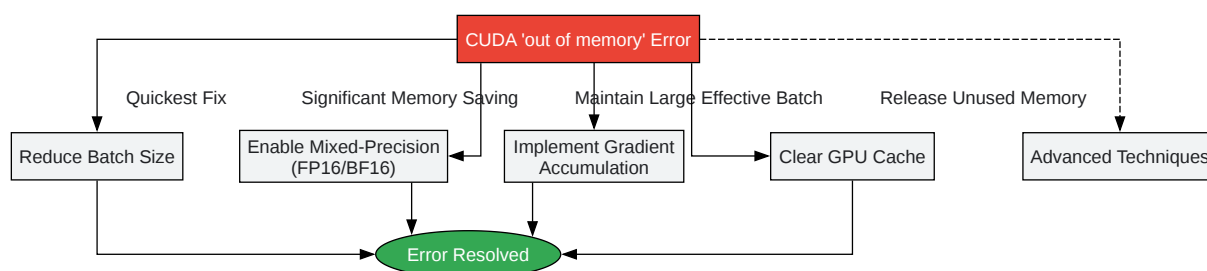
This is one of the most common errors encountered when working with large models and datasets on **H100** GPUs.

1. What are the immediate steps to take when I encounter a "CUDA out of memory" error?

When you encounter this error, it means your program tried to allocate more memory on the GPU than is available. Here's a systematic approach to troubleshoot this issue:

- **Reduce Batch Size:** This is often the quickest fix. A smaller batch size consumes less memory.^{[1][2][3]} However, be aware that this might affect model convergence, so you may need to adjust the learning rate.

- Enable Mixed-Precision Training: Utilize the **H100**'s Tensor Cores by using lower precision data types like FP16 or BF16.[4][5][6][7][8][9] This can halve the memory footprint of your model and data.
- Use Gradient Accumulation: This technique allows you to simulate a larger batch size by accumulating gradients over several smaller batches before performing a weight update.[10][11][12][13][14]
- Clear GPU Memory Cache: In frameworks like PyTorch, explicitly clear the cache to free up unused memory.[15][16]



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

Caption: Initial steps for resolving "CUDA out of memory" errors.

2. How do I implement mixed-precision training in my workflow?

Mixed-precision training combines the use of 16-bit and 32-bit floating-point types to reduce memory usage and accelerate computation, often without a significant loss in model accuracy. [9]

Experimental Protocol for Mixed-Precision Training:

- Framework Support: Utilize built-in automatic mixed precision (AMP) functionalities in your deep learning framework.

- PyTorch: Use `torch.cuda.amp`.[\[6\]](#)[\[7\]](#)[\[8\]](#)[\[17\]](#)
- TensorFlow: Use `tf.keras.mixed_precision.set_global_policy('mixed_float16')`.[\[6\]](#)[\[7\]](#)[\[18\]](#)
- Gradient Scaling: To prevent underflow of small gradient values during backpropagation with FP16, use a loss scaler. AMP in PyTorch and TensorFlow handles this automatically.
- Monitor for NaNs: Keep an eye on model outputs and loss values for "Not a Number" (NaN) results, which can indicate numerical instability. If this occurs, you may need to keep certain sensitive layers in FP32.

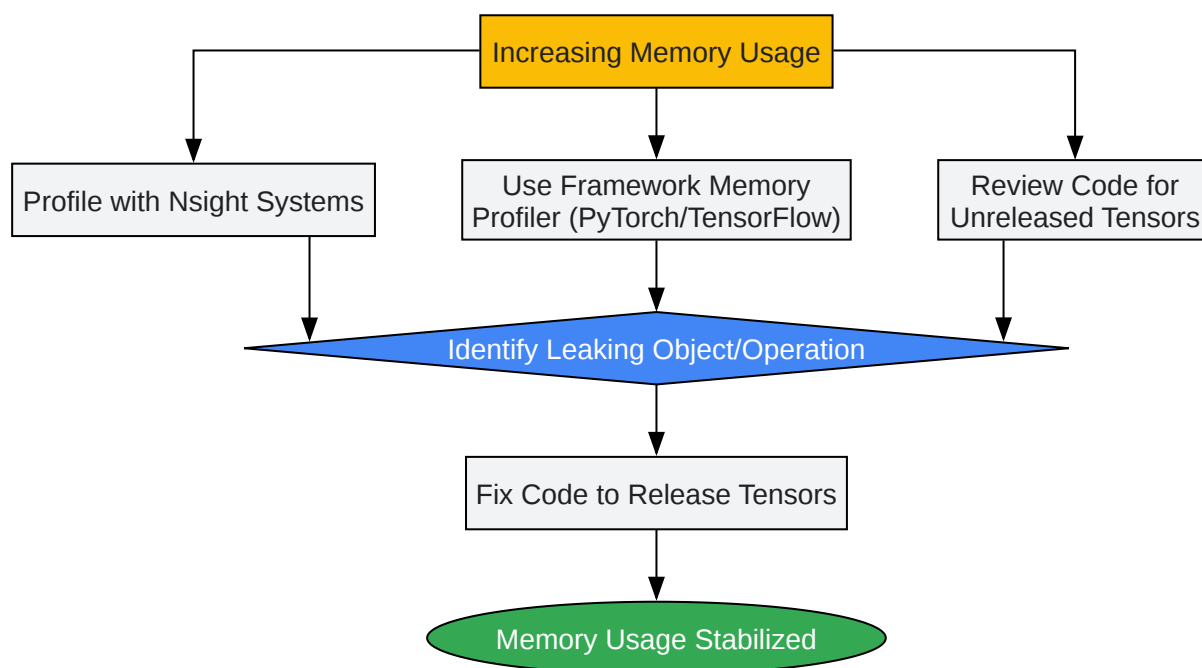
Precision	Memory Usage (per parameter)	H100 Support	Typical Use Case
FP64	8 bytes	Yes	High-precision scientific computing
FP32	4 bytes	Yes	Standard for many deep learning models
TF32	4 bytes (uses 19 bits)	Yes (Tensor Cores)	Default for dense and convolutional layers
BF16	2 bytes	Yes (Tensor Cores)	Good for training, less prone to underflow than FP16
FP16	2 bytes	Yes (Tensor Cores)	Reduces memory by half compared to FP32 [9]
FP8	1 byte	Yes (Transformer Engine)	Accelerates transformer models [6] [7]

Issue: Gradual Increase in Memory Usage Leading to a Crash

If you observe your application's memory consumption steadily growing over time, you might be facing a memory leak or fragmentation issue.

1. How can I identify the cause of increasing memory usage?

- Profiling Tools: Use NVIDIA's Nsight Systems to profile your application and visualize memory allocation and deallocation patterns.[\[4\]](#)[\[7\]](#)[\[19\]](#)
- Framework-Specific Tools:
 - PyTorch: Use `torch.cuda.memory_summary()` to get a detailed report of memory usage.[\[20\]](#)
 - TensorFlow: Utilize the TensorFlow Profiler to trace memory usage over time.
- Code Review: Look for tensors that are not being released, especially in long-running loops.[\[2\]](#)[\[21\]](#) Ensure that any custom CUDA kernels properly deallocate memory.[\[22\]](#)



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

Caption: Workflow for diagnosing memory leaks.

2. What is memory fragmentation and how can I mitigate it?

Memory fragmentation occurs when memory is allocated and deallocated, leaving small, non-contiguous blocks of free memory.^{[1][2][22]} Even if the total free memory is sufficient, a large contiguous block may not be available for a new allocation request.

Mitigation Strategies:

- **Memory Pooling:** Deep learning frameworks like PyTorch and TensorFlow use caching allocators that manage a pool of memory to reduce fragmentation.^{[15][23]}
- **Pre-allocation:** For applications with predictable memory requirements, pre-allocating a large chunk of memory at the beginning can help.
- **Restart the Process:** In some cases, the simplest solution for a long-running process that has become fragmented is to periodically save the state and restart.

Frequently Asked Questions (FAQs)

Q1: What are the most common causes of memory errors on **H100** GPUs?

The most frequent causes include:

- **Insufficient GPU Memory:** The model or data is simply too large for the **H100**'s 80GB of HBM3 memory.^{[1][22]}
- **Large Batch Sizes:** Using a batch size that exceeds the available memory.^{[1][2]}
- **High-Precision Data Types:** Using FP32 or FP64 when lower precisions would suffice.^[1]
- **Memory Fragmentation:** Continuous allocation and deallocation of memory leading to unusable free memory blocks.^{[1][22]}

- Software and Framework Overhead: Deep learning frameworks allocate extra memory for caching and intermediate computations.[\[1\]](#)

Q2: How can I optimize my data loading pipeline to reduce memory usage?

- Data Streaming: Avoid loading the entire dataset into GPU memory at once. Stream batches from the CPU as needed.[\[1\]](#)
- Pin Memory: In PyTorch's DataLoader, set `pin_memory=True` to speed up data transfers from CPU to GPU.[\[8\]](#)[\[17\]](#)
- NVIDIA DALI: Use the NVIDIA Data Loading Library (DALI) to offload data preprocessing to the GPU, which can be more memory-efficient.[\[8\]](#)

Q3: What are some advanced techniques to fit very large models on an **H100** GPU?

- Gradient Checkpointing (Activation Recomputation): This technique trades compute for memory by not storing intermediate activations in the forward pass and recomputing them during the backward pass.[\[4\]](#)[\[5\]](#)[\[17\]](#)[\[21\]](#)
- Model Parallelism: For extremely large models, you can split the model itself across multiple GPUs. The 4th Gen NVLink on the **H100** is beneficial for this.[\[4\]](#)
- Memory-Efficient Optimizers: Some optimizers, like AdamW, can be more memory-efficient than others.[\[4\]](#)
- FlashAttention: For transformer models, use memory-efficient attention mechanisms like FlashAttention.[\[21\]](#)[\[24\]](#)

Q4: How do I choose the right precision for my application?

Precision	When to Use	Considerations
FP64	When high numerical precision is critical (e.g., certain scientific simulations).	Doubles memory usage compared to FP32.
FP32	As a baseline for model training and when numerical stability is a concern.	Standard precision, but can be memory-intensive.
TF32	Default for many operations on H100. Good balance of performance and precision.	Not a true 32-bit float, but generally provides similar accuracy.
BF16	Recommended for deep learning training due to its dynamic range.	Better stability than FP16 for training.
FP16	For inference and training where memory is highly constrained.	Can suffer from underflow/overflow; requires loss scaling.
FP8	For accelerating transformer models, especially during inference. ^{[6][7]}	Requires the Transformer Engine on the H100.

Q5: Where can I find more information and tools for debugging **H100** memory issues?

- NVIDIA Developer Documentation: The official CUDA and **H100** documentation provides in-depth information.
- NVIDIA Nsight Tools: A suite of tools for profiling and debugging GPU applications.^{[4][7]}
- Framework-Specific Documentation: The PyTorch and TensorFlow documentation have detailed sections on CUDA memory management.^{[15][25]}

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. What are the common causes of out-of-memory errors on A100 and H100 GPUs? - Massed Compute [massedcompute.com]
- 2. What are the common causes of CUDA out of memory errors - Massed Compute [massedcompute.com]
- 3. saturncloud.io [saturncloud.io]
- 4. What are the recommended memory allocation strategies for mixed-precision training on H100 GPUs? - Massed Compute [massedcompute.com]
- 5. How do I optimize the memory usage of the NVIDIA H100 PCIe GPU when running large TensorFlow models? - Massed Compute [massedcompute.com]
- 6. Optimizing Deep Learning Pipelines with NVIDIA H100 [centron.de]
- 7. Optimizing deep learning pipelines for maximum efficiency | DigitalOcean [digitalocean.com]
- 8. How do I optimize my PyTorch model for the NVIDIA H100 PCIe GPU? - Massed Compute [massedcompute.com]
- 9. Train With Mixed Precision - NVIDIA Docs [docs.nvidia.com]
- 10. Day 24: Gradient Accumulation - DEV Community [dev.to]
- 11. medium.com [medium.com]
- 12. Gradient Accumulation: Overcome GPU Memory Limitations | Vultr Docs [docs.vultr.com]
- 13. Save Memory on GPU — mmengine 0.11.0rc0 documentation [mmengine.readthedocs.io]
- 14. medium.com [medium.com]
- 15. CUDA semantics — PyTorch 2.9 documentation [docs.pytorch.org]
- 16. medium.com [medium.com]
- 17. Can you provide examples of memory allocation strategies for A100 and H100 GPUs in PyTorch and TensorFlow? - Massed Compute [massedcompute.com]
- 18. How do I optimize the memory configuration of my A100 or H100 GPU for my TensorFlow application? - Massed Compute [massedcompute.com]
- 19. How do I configure my system to allocate more memory to the NVIDIA A100 and H100 GPUs? - Massed Compute [massedcompute.com]

- 20. PyTorch 101 Memory Management and Using Multiple GPUs | DigitalOcean [digitalocean.com]
- 21. How do I troubleshoot NVIDIA H100 GPU memory issues caused by software bugs in the large language model training code? - Massed Compute [massedcompute.com]
- 22. What are the common causes of memory-related issues on the H100? - Massed Compute [massedcompute.com]
- 23. What are the optimal memory allocation settings for NVIDIA H100 GPUs in deep learning applications? - Massed Compute [massedcompute.com]
- 24. VRAM in Large Language Models: Optimizing with NVIDIA H100 VRAM GPUs [uvation.com]
- 25. apxml.com [apxml.com]
- To cite this document: BenchChem. [H100 GPU Technical Support Center: Troubleshooting Memory Allocation Errors]. BenchChem, [2025]. [Online PDF]. Available at: [https://www.benchchem.com/product/b15585327#troubleshooting-memory-allocation-errors-on-h100-gpus]

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

