

---

# Week 3 Reading Summaries: GPU Performance + MI300X vs H100 Readings

---

**Amogh Patankar**

University of California, San Diego  
3869 Miramar Street Box #3037, La Jolla, CA- 92092  
apatankar@ucsd.edu

## Abstract

1       The articles for this week’s reading summary describe the background on graphical  
2       processing units (GPUs), which are common modern hardware to execute machine  
3       learning workloads efficiently. Moreover, we read about operation execution and  
4       limitations, and matrix multiplication and its performance on GPUs and in other ML  
5       models. Finally, we also read about benchmarking of various hardware platforms,  
6       namely, H100, MI300X, and H200, with respect to latency, throughput, and total  
7       cost of ownership (TCO).

## 8    **1    TensorFlow**

### 9    **1.1   Introduction**

10   Graphical processing units are the basic building blocks of modern ML hardware, and there needs  
11   to be a proper understanding of what a GPU entails. Namely, we look at the structure/architecture,  
12   operation execution, performance estimation and limitations, and applications of deep learning and  
13   their respective limitations when utilizing GPUs.

### 14   **1.2   GPU Architecture Fundamentals**

15   The GPU has a parallel processing architecture, with processing and memory hierarchies. To  
16   summarize, it consists of streaming multiprocessors, on chip L2 cache, and DRAM. The SMs,  
17   DRAM, and L2 are responsible for arithmetic and instructions, data, and code access, respectively.

18   SMs have schedulers and instruction execution pipelines, with MAC/MAD operation being the most  
19   frequent in modern ML; data types have different MAC ops/clock, with a MAC consisting of two  
20   operations. So, to calculate FLOPs, we multiple throughput by two, and for GPU FLOPs, we multiple  
21   that value by the # of SMs and SM clock rate.

22   Nvidia GPUs allow various datatype execution, like INT8 and FP16; tensor cores are used to speed  
23   up matmul, MAC, and other operations on small matrices (eg. 4x4). In doing so, the computation can  
24   occur in higher precision, and non-executable ops are executed in other cores.

### 25   **1.3   GPU Execution Model**

26   There are two main concepts to connect thread count to GPU performance: 1) grouping threads to  
27   equally sized thread blocks, and execution [thread] switching. These two put together make up the  
28   nuances of the execution model of GPUs.

29   Because GPUs contain many SMs, each has pipelines to execute threads, and in turn, have shared  
30   memory. This results in a two-level thread hierarchy. To fully utilize a GPU when executing a

31 function, we'd use multiple thread blocks, and the number of threads should be a multiple of the # of  
32 SMs.

33 To maximize efficiency, we launch functions executing in waves of thread blocks, minimizing tail  
34 effect (underutilization of GPU).

## 35 1.4 Understanding Performance

36 GPU Performance is limited by memory bandwidth, math bandwidth, and latency. Memory time = (#  
37 of bytes accessed in memory)/(process memory bandwidth). Math time = (# of ops)/(processor math  
38 bandwidth). All in all, algos are math limited if arithmetic intensity > processor  $BW_{math}/BW_{mem}$ .

## 39 1.5 DNN Operation Categories

### 40 1.5.1 Elementwise Ops

41 These are unary/binary ops, layers in this category perform mathematical ops independent of other  
42 tensor values.

### 43 1.5.2 Reduction Ops

44 These ops produce values computed over a range of input tensor vals. Examples include pooling  
45 layers, and SoftMax.

### 46 1.5.3 Dot-Product Ops

47 These ops are expressed as dot-products of elements from two tensors. Examples include FC layers,  
48 convolutions, and matmuls.

## 49 1.6 Summary

50 When evaluating GPU performance, we look at following criteria:

- 51 • # of SMs on GPU, determine ops:bytes ratio (i.e.  $\frac{BW_{math}}{BW_{mem}}$ )
- 52 • Compute arithmetic intensity, check parallelism with thread blocks and SM count.

53 Likely performance limiters are often latency, math, and/or memory (last two depend on ops:byte  
54 ratio).

## 55 2 Matrix Multiplication Background User Guide

### 56 2.1 Background

57 General Matrix Multiplications, known as GEMMs, make up neural network operations, in FC,  
58 recurrent, and convolutional layers alike. GEMMs are essentially an operation defined as  $C = \alpha AB +$   
59  $\beta C$ , including two matrices, two scalar inputs, and a pre-existing matrix.

### 60 2.2 Math and Memory Bounds

61 Given a product of two matrices, we have a total of  $M * N * K$  fused MACs, i.e.  $2 * M * N * K$ , as  
62 each MAC is 2 FLOPs. The arithmetic intensity of a GEMM is as follows:  $\frac{2 * M * N * K}{2(M * K + N * K + M * N)}$   
63

#### 64 2.2.1 GPU Implementation

65 GEMMs' output matrices are partitioned into tiles, assigned to thread blocks on GPU. We step  
66 through K dimension in tiles, loading values from A and B, and performing MAC.

## 67 2.2.2 Tensor Core Requirements

68 Tensor cores maximize speed of tensor multiplies, and performance of GEMMs improves when  
69 M, N, and K are aligned to multiples of 16 bytes. For example, with FP16 data, each element is 2  
70 bytes, hence dimensions should be multiples of 8. Scaling down (multiples of 16, then 8, then 4) is  
71 recommended.

## 72 2.2.3 Typical Tile Dimensions in cuBLAS and Performance

73 Larger tiles reuse more data, using less bandwidth, and increase efficiency compared to smaller tiles.  
74 However, you run fewer large tiles in parallel, which may lead to tailing problem as mentioned earlier.  
75 Tradeoff between tile efficiency and parallelism suggests that larger GEMMs impact the tradeoff less.  
76 Flipped, a smaller GEMM causes a larger tradeoff in the sense that the GPU won't run at maximum  
77 efficiency.

## 78 2.3 Dimension Quantization Effects

79 Tile and wave quantization affect execution efficiency.

### 80 2.3.1 Tile Quantization

81 This is when matrix dimensions aren't divisible by thread block tile size; for example, with 256x257  
82 data vs 256x256 data using 128x128 tiles, utilization will be 6 vs 4 tiles, a 50% increase in ops, when  
83 in reality only has .39% more data.

### 84 2.3.2 Wave Quantization

85 Tile quantization quantizes problem size to the size of each tile, wave quantization quantizes the total  
86 # of tiles to the # of multiprocessors on the GPU. Essentially, wave quantization quantizes the work  
87 to the size of the GPU, not the size of tiles of thread blocks on the GPU.

## 88 3 MI300X vs H100

### 89 3.1 Intro

90 This article compares the performance of AMD and Nvidia GPU processors, in terms of specs and  
91 TCO. They find that the advantage of the MI300X isn't valid due to a lack of AMD's SW stack and  
92 less testing, including bugs.

### 93 3.2 Key Findings

94 They provide many key findings, but essentially, it all boils down to this: Nvidia has an incredible  
95 experience out-of-the-box (OOTB), while AMD struggles mightily. They compared TCO and  
96 performance, where AMD has a lower TCO but has worse performance too. They also note that  
97 AMD requires better coverage with regards to ROCm, internal testing, debugging, benchmarking,  
98 and scalability.

### 99 3.3 AMD vs Nvidia

100 While the H100 was launched in late 2022, AMD's alternative, MI300X was introduced in late 2023.  
101 It has a extremely low TCO on paper, and represents higher performance as well (on paper).

### 102 3.4 General Matrix Multiply (GEMM)

103 When benchmarking for GEMM for BF16, the H100 and 200 achieve 720 TFLOP/s, while the  
104 MI300X hits only 620 TFLOP/s. When compared to their marketed values, Nvidia struggles by only  
105 about 20%, while AMD drops by 50%. This gap only widens with FP8, looking at a 20% gap as  
106 opposed to 10% between Nvidia and AMD.

### 107 **3.5 HBM Performance**

108 The MI300X has better memory bandwidth compared to the H100/200, useful for inferencing and  
109 sometimes training.

### 110 **3.6 Training + Testing Methodology**

111 SemiAnalysis used a unique strategy to test performance; they used MLPerf’s GPT3 175B Training to  
112 measure train time to convergence, but was difficult to run. Hence, they ran GPT 1.5B DDP, Llama3  
113 8B, Llama3 70B, and Mistral 7B v0.1.

### 114 **3.7 Single Node Training Performance**

115 When comparing single node training performance, SemiAnalysis found the following:

116

- 117 • BF16: MI300X underperforms H100/200 for GPT/Llama3-8B/Mistral 7B, barely outper-  
118 forms Nvidia on Llama3-70B using a private repo (not valid).
- 119 • FP8: MI300X underperforms on all benchmarks

### 120 **3.8 Multi-Node Training Performance**

121 Regardless of precision, H100 crushes the MI300X, and the H200 wasn’t available for multinode  
122 training.

### 123 **3.9 Scale Up NVLink/xGMI Topology**

124 Scaleup fabric is important for GPU clusters as it provides fast tensor parallelism. As such, NVLink  
125 for Nvidia provides 450GB/s per GPU with 8 GPUs, while xGMI provides 448GB/s for 8 GPUs for  
126 AMD. However, xGMI provides point-to-point bandwidth, so for 8 GPUs, we only get 64 GB/s per  
127 per GPU. For NVLink, it’s a switched topography, meaning we get full 450 GB/s due to NVSwitches.

### 128 **3.10 Single Node NCCL Collective**

129 For all tests undertaken, the H100 beats the MI300X, regardless of # of GU, precision, etc.

### 130 **3.11 MultiNode RCCL/NCCL Collectives + Scaleout Network Benchmarks**

131 For all the GPUs here (H100/200 + MI300X), each GPU connects to other nodes using scaleout  
132 network, with a 400G NIC. In any case, Nvidia’s GPU’s performance beats MI300X, due to AMD’s  
133 RCCL being a fork of Nvidia’s NCCL, restricting access to up-to-date features.

### 134 **3.12 AMD Forked Libraries**

135 SemiAnalysis points out that AMD forks off of Nvidia’s ecosystem libraries. For NCCL, APex,  
136 Megatron, and Pytorch Eagers Kernel, AMD forks to ROCm/RCCL, Apex, Megatron, and PyTorch  
137 Eagers Kernel. This is a huge limitation for AMD MI300X performance due to SW restrictions.

### 138 **3.13 Recommendations to fix AMD SW**

139 SemiAnalysis provides a thorough recommendation to AMD in order to fix their software stack to  
140 allow for optimal MI300X performance. To summarize, they conclude that SW engineers (#, quality),  
141 lack of resources available to AMD engineers, integration of ROCm with Meta’s PyTorch team, and  
142 using MLPerf/TorchBench to benchmark.