

Parallelize Muon with FSDP2

Motivation

Algorithm 1 Distributed Muon

Require: Full Gradients G , DP partitioned Momentum m , DP partitioned parameters p , momentum μ .

```

1: // Reduce-scatter  $G$  on DP for correct gradients
2:  $g = \text{reduce\_scatter}(G, \text{dp\_group})$ 
3: // Apply momentum to  $g$  using local partitioned momentum  $m$ 
4:  $g' = \text{update\_with\_momentum}(g, m, \mu)$ 
5: // DP Gather: gathering  $g'$  across DP into a full matrix  $G$ 
6:  $G = \text{gather}(g', \text{dp\_group})$ 
7: // Calculate Muon update
8:  $U = \text{Newton-Schulz}(G)$ 
9: // Discard the rest of  $U$  and only keep the local partition  $u$ , then apply the update rule
10:  $p' = \text{apply\_update}(p, u)$ 
11: // All-gather updated  $p'$  into  $P$ 
12:  $P = \text{all\_gather}(p', \text{dp\_group})$ 
13: // Return the update RMS for logging
14: return  $\sqrt{u^2.\text{mean}()}$ 

```

Distributed Muon by Moonlight

While a distributed version of Muon is available, it has the drawback of redundant computations across GPUs.

GPU 0		C[0]	C[1]	C[2]	C[3]	
GPU 1		C[0]	C[1]	C[2]	C[3]	
GPU 2		C[0]	C[1]	C[2]	C[3]	
COMM	AG[0]		AG[1]		AG[2]	
					AG[3]	
						AG[4]

Execution timeline of Distributed Muon

- $C[i]$: Compute Newton-Schulz(G) for i -th gradient
- $AG[i]$: AllGather i -th gradient
- $G[i]$: Gather i -th gradient
- $SC[i]$: Scatter i -th gradient

Algorithm

Parallel Muon

Algorithm 1 Parallel Muon

Require: DP partitioned gradient \mathbf{g} , DP partitioned Momentum \mathbf{m} , DP partitioned parameter \mathbf{p} , momentum μ , local rank r

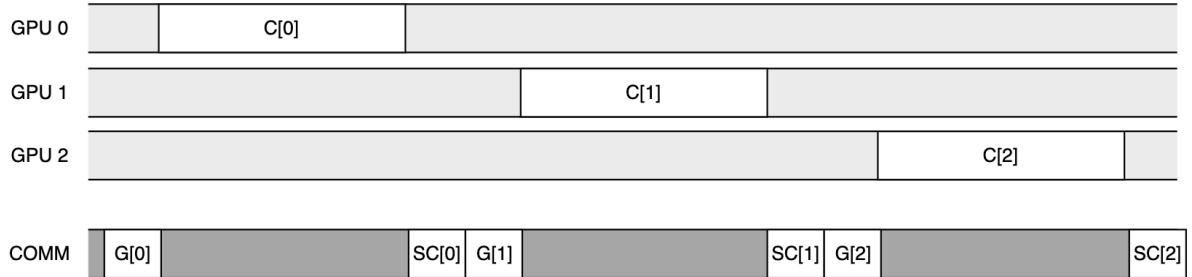
```

1: // Apply momentum to  $\mathbf{g}$  using local partitioned momentum  $\mathbf{m}$ 
2:  $\mathbf{g}' \leftarrow \text{update\_with\_momentum}(\mathbf{g}, \mathbf{m}, \mu)$ 
3: // Schedule  $\mathbf{g}'$  to rank  $\mathbf{R}$ 
4:  $\mathbf{R} \leftarrow \text{schedule}(\mathbf{g}', \text{dp\_group})$ 
5: // Gather  $\mathbf{g}'$  across DP into a full matrix  $\mathbf{G}$  to rank  $\mathbf{R}$ 
6:  $\mathbf{G} \leftarrow \text{gather}(\mathbf{g}', \text{dp\_group}, \text{dst}=\mathbf{R})$ 
7: // Calculate Newton-Schulz only in  $\mathbf{R}$ 
8: if  $r == \mathbf{R}$  then
9:    $\mathbf{u} \leftarrow \text{Newton-Schulz}(\mathbf{G})$ 
10: else
11:    $\mathbf{u} \leftarrow \text{None}$ 
12: end if
13: // Scatter a full matrix  $\mathbf{u}$  across DP
14:  $\mathbf{u}' \leftarrow \text{scatter}(\mathbf{u}, \text{dp\_group}, \text{src}=\mathbf{R})$ 
15: // Apply DP partitioned  $\mathbf{u}'$  to  $\mathbf{p}$ 
16:  $\mathbf{p}' \leftarrow \text{apply\_update}(\mathbf{p}, \mathbf{u}')$ 
17: return  $\mathbf{p}'$ 

```

We eliminate redundant computation by assigning each parameter to a specific GPU.

However, without proper scheduling, this optimization can lead to poor GPU utilization. In particular, although redundant computation is avoided by assigning each parameter to a specific rank, it causes idle time—since all other ranks must wait for the scatter communication to complete before proceeding.

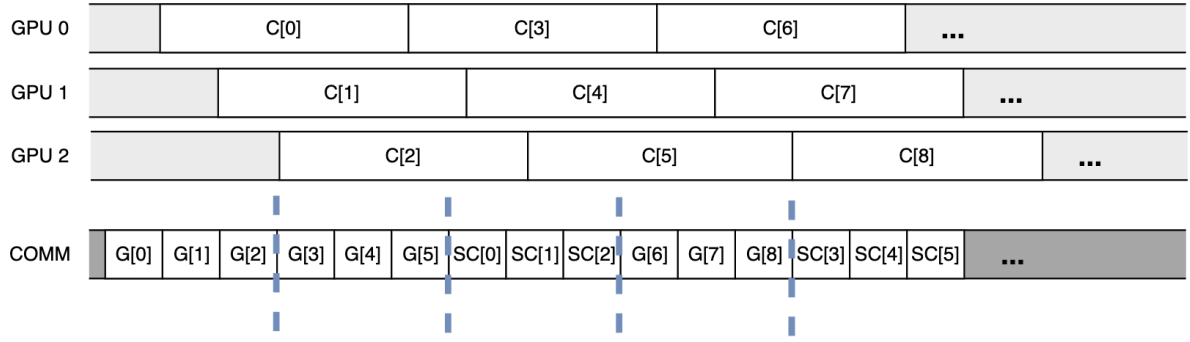


Execution timeline of Parallel Muon

Scheduling Sub-Operations

We can schedule the whole sub-operations as follows, due to the following reasons:

- There are no dependencies between parameters.
- GPUs can execute computation and communication concurrently.



Execution timeline of re-scheduled Parallel Muon

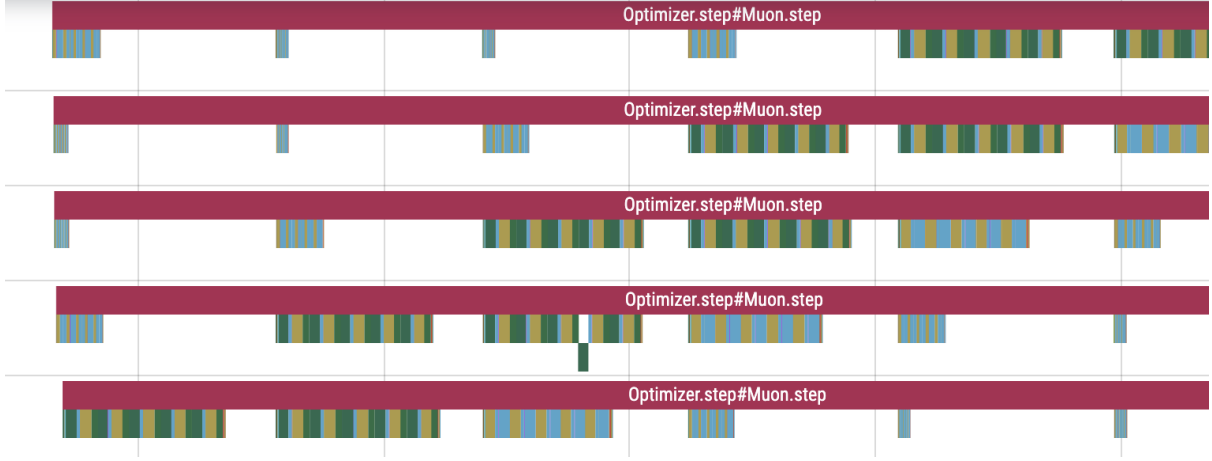
We define the chunk size C as the number of GPUs and schedule each sub-operation in batches of size C . This scheduling allows each GPU to continue computation even while waiting for collective communication to complete.

[Algorithm] (To be written)

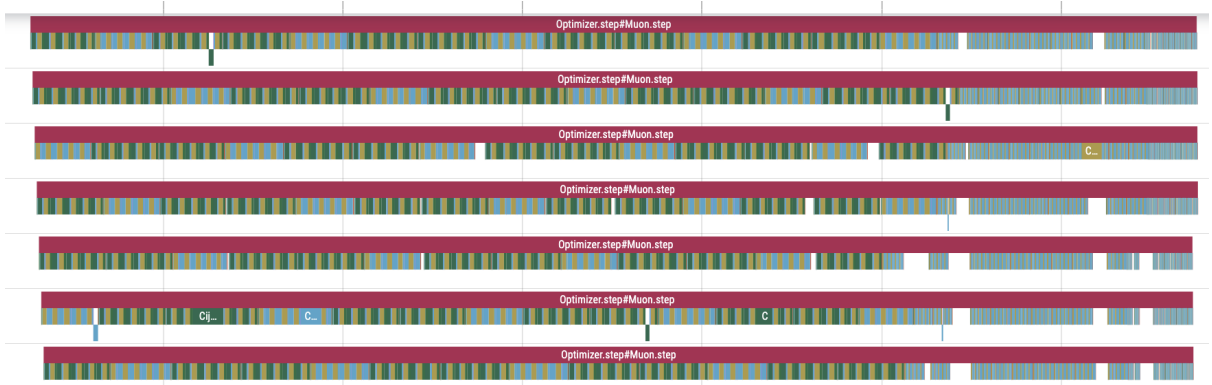
Load Balancing

If parameters in a chunk have imbalanced computation loads, idle bubbles may occur. To mitigate this, we apply load balancing based on per-parameter FLOPs.

Imbalanced (Round Robin)



After Load Balancing



Implementation

The full implementation is available in `optimizer/torch-ext/optimizer/muon.py`. To enable concurrent computation and communication, we use separate compute and communication streams (`torch.cuda.Stream`) and use `torch.cuda.Event` to synchronize between sub-operations.

Thanks to the simplicity of `torch.DTensor` and `torch.distributed`, the implementation remains straightforward and low in complexity.

Evaluation

We evaluated the performance using 10B model currently in development, achieving 151 TFLOPS per GPU during the optimizer step.

Model Size	TFLOPs for Muon	GPUs	Elapsed time	TFLOPS/GPU
10B	847.45	4xMI250 (8 devices)	1.4 s	151

Based on the breakdown, 7% of the time is attributed to updating sharded gradients and parameters, 78% to GEMM operations, and the remaining 15% to non-overlapped communication overhead.