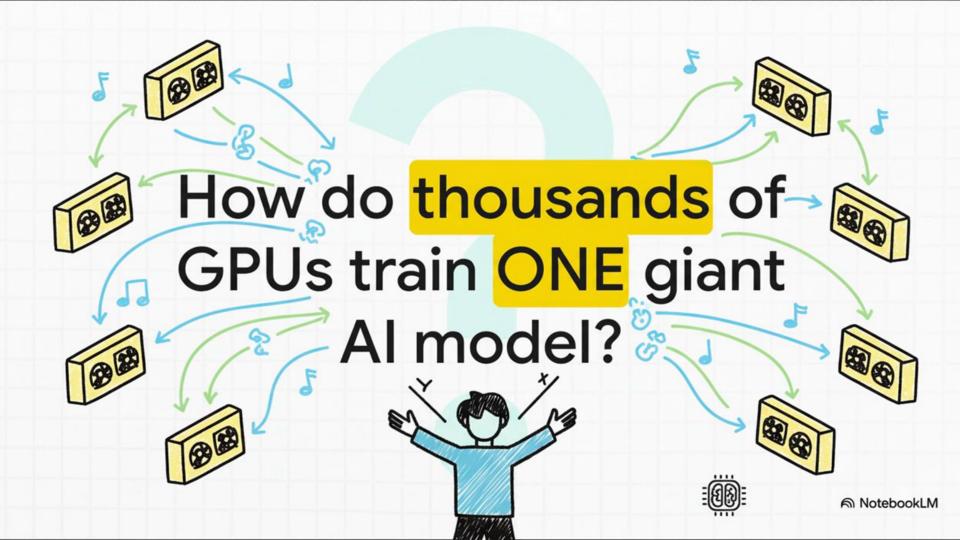
Demystifying NCCL: An In-depth Analysis of GPU Communication Protocols and Algorithms

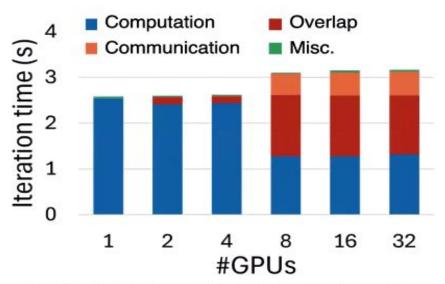
Authors: Zhiyi Hu, Siyuan Shen, Tommaso Bonato, Sylvain Jeaugey, Cedell Alexander, Eric Spada, James Dinan, Jeff Hammond, Torsten Hoefler





GPU Communication

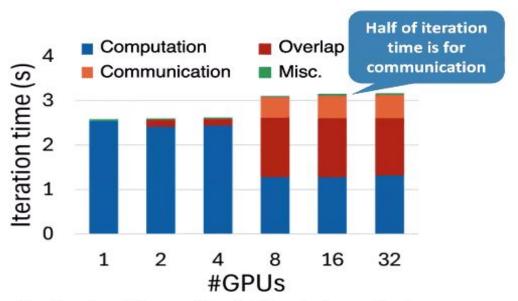
GPU Communication Bottleneck



Iteration time breakdown of the full fine-tuning method on a 1B parameters Transformer decoder model [1]

[1] N. Alnaasan et al., "Characterizing Communication in Distributed Parameter-Efficient Fine-Tuning for Large Language Models," HOTI 2024, pp. 11–19.

GPU Communication Bottleneck



Iteration time breakdown of the full fine-tuning method on a 1B parameters Transformer decoder model [1]

[1] N. Alnaasan et al., "Characterizing Communication in Distributed Parameter-Efficient Fine-Tuning for Large Language Models," HOTI 2024, pp. 11–19.

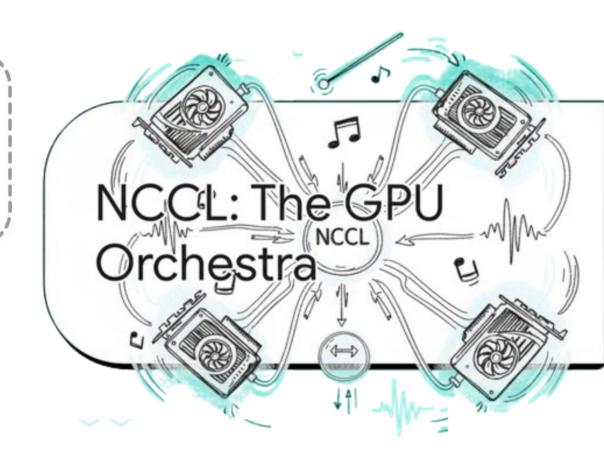




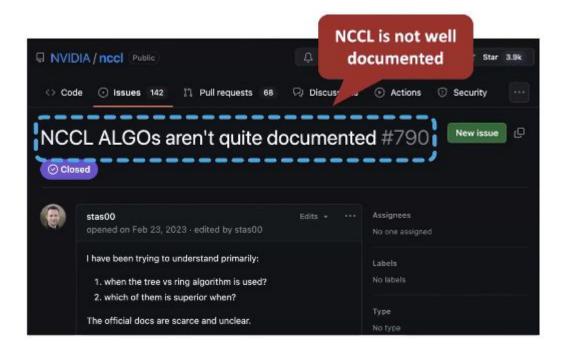
TensorFlow

NCCL

NVIDIA GPU



Motivation



Communicator Creation

```
// Create communicator
ncclComm_t comm;
ncclCommInitRank(&comm, nranks, id, rank);
```

Communicator Creation

```
// Create communicator
ncclComm_t comm;
ncclCommInitRank(&comm, nranks, id, rank);
```

Collective Communication / P2P Communication

```
// Collective communication call
ncclAllReduce(sendbuff, recvbuff, count, ncclFloat, ncclSum, comm, stream);
// Point-to-point communication call
ncclSend(sendbuff, count, ncclFloat, next_rank, comm, stream);
ncclRecv(recvbuff, count, ncclFloat, prev_rank, comm, stream);
```



```
// Create communicator
ncclComm t comm;
ncclCommInitRank(&comm, nranks, id, rank);
// Start group operation
ncclGroupStart();
// Collective communication call
ncclAllReduce(sendbuff, recvbuff, count, ncclFloat, ncclSum, comm, stream);
// Point-to-point communication call
ncclSend(sendbuff, count, ncclFloat, next_rank, comm, stream);
ncclRecv(recvbuff, count, ncclFloat, prev_rank, comm, stream);
// End group operation
ncclGroupEnd();
```

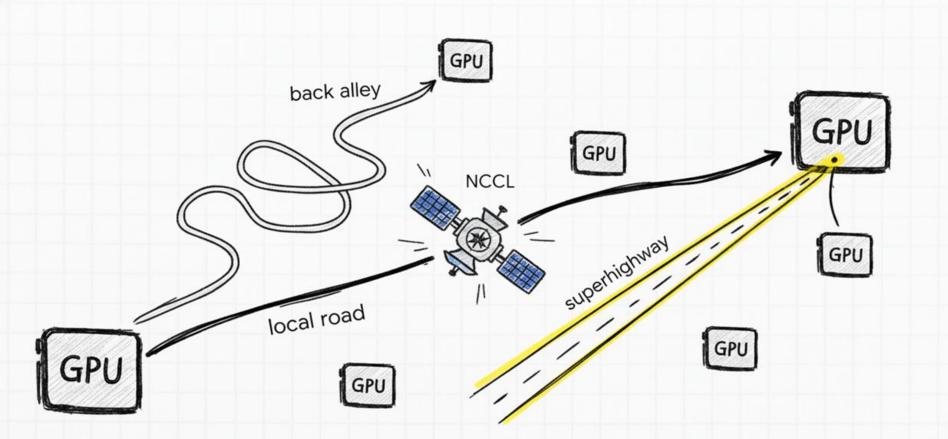


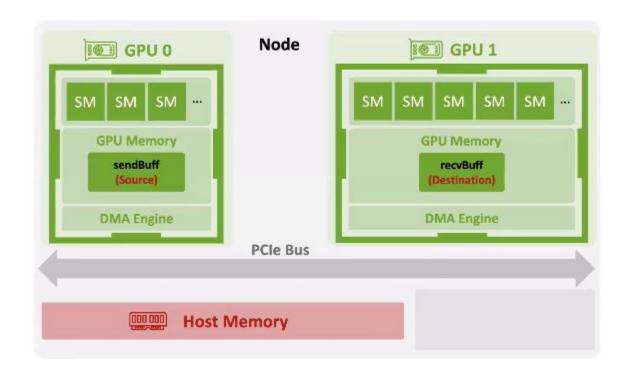
```
// Create communicator
ncclComm_t comm;
ncclCommInitRank(&comm, nranks, id, rank);
// Start group operation
ncclGroupStart();
// Collective communication call
ncclAllReduce(sendbuff, recvbuff, count, ncclFloat, ncclSum, comm, stream);
// Point-to-point communication call
ncclSend(sendbuff, count, ncclFloat, next_rank, comm, stream);
ncclRecv(recvbuff, count, ncclFloat, prev rank, comm, stream);
// End group operation
ncclGroupEnd();
// Destroy communicator
ncclCommDestroy(comm);
```

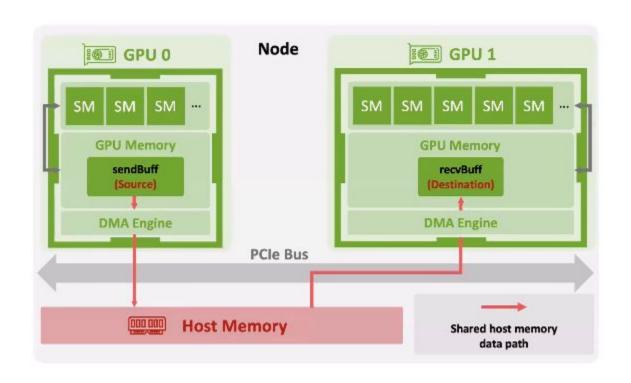
Communication Protocols

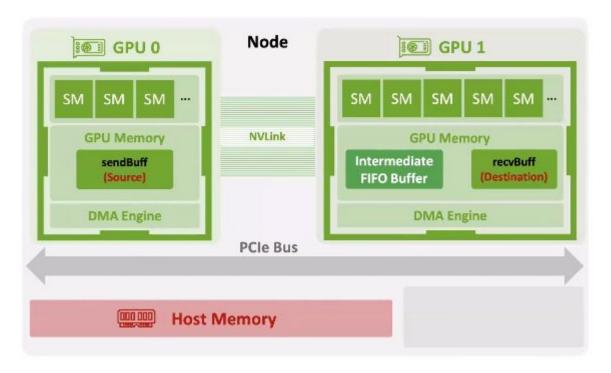
Protocol	Design Goal	Payload	Bandwidth Utilization	Latency Per-Hop
Simple	High Bandwidth	Data Chunks	Near Peak	~6µs
Low Latency (LL)	Low Latency	4B data + 4 B flag	25~50% of peak	~1µs
LL128	Best of both	120B data + 8B flag	~95% of peak	~2µs

Note: LL128 requires NVLink

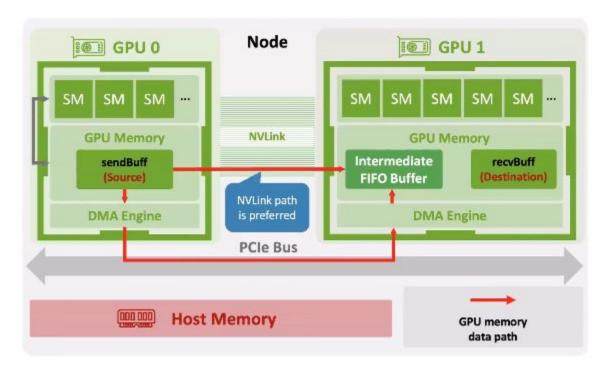




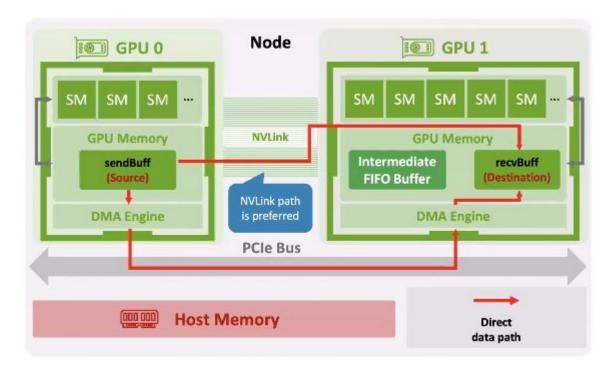




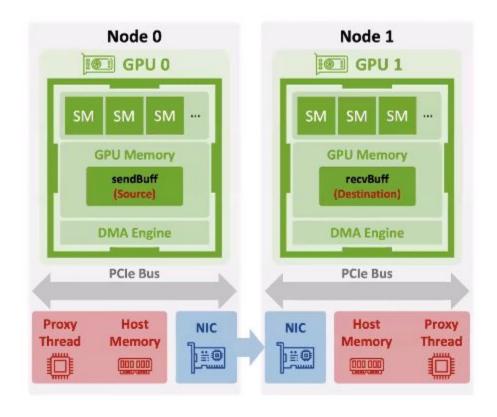
GPU Memory

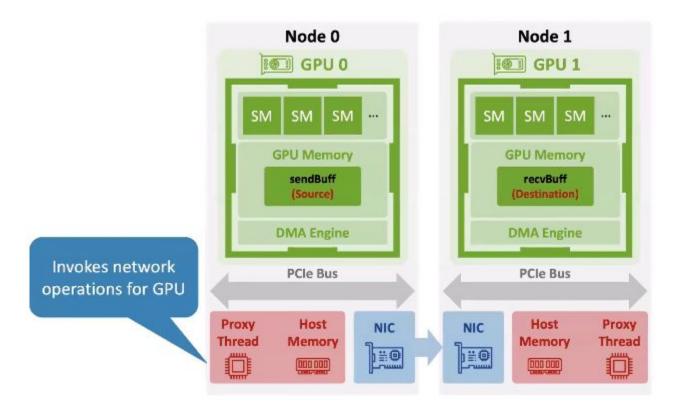


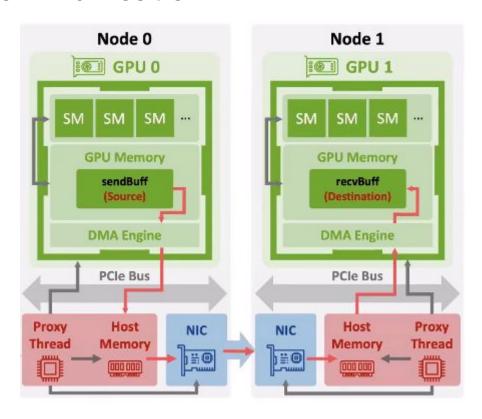
GPU Memory

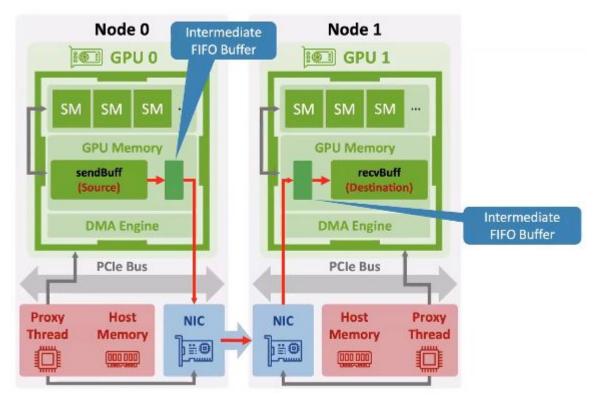


P2P_DIRECT mode enabled and GPUs belong to the same process



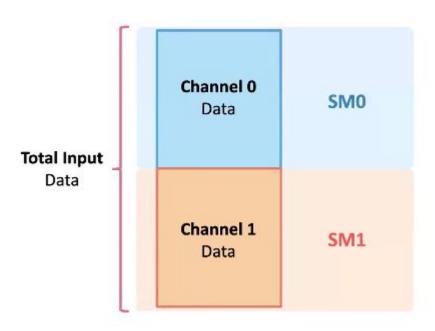




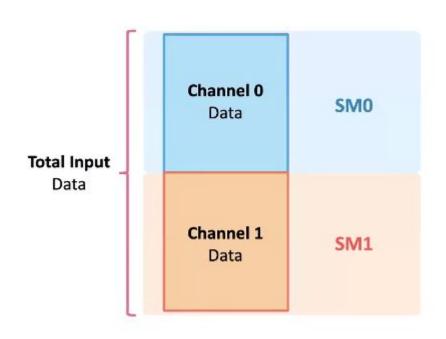


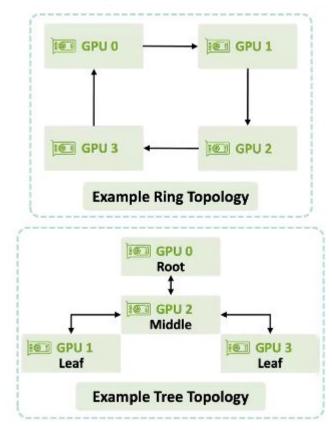
With GPU Direct RDMA

Communication Channels & Logical Topology

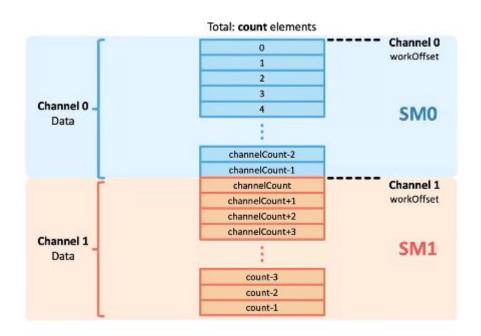


Communication Channels & Logical Topology

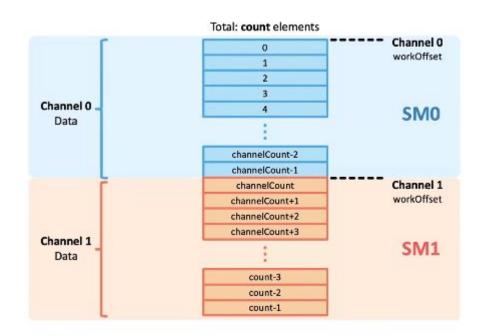




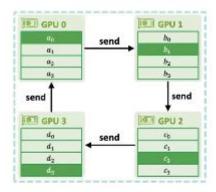
Iterative Execution and Communication Primitives

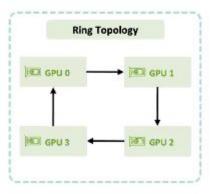


Iterative Execution and Communication Primitives

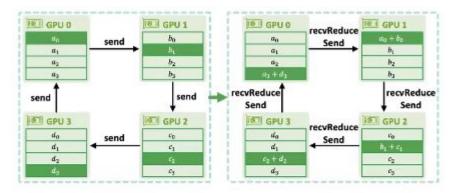


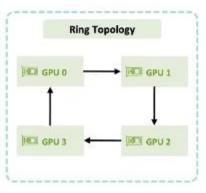
Step Index	NCCL Primitive	
0	send	
1 to $k-2$	recvReduceSend	
k-1	recvReduceCopySend	
k to 2k-3	recvCopySend	
2k-2	recv	





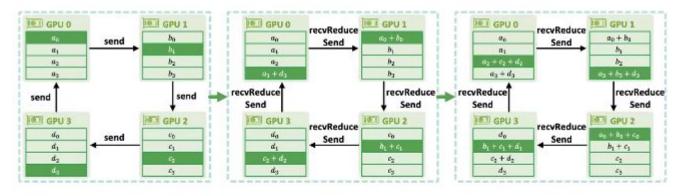
Step Index	NCCL Primitive	
0	send	
1 to $k-2$	recvReduceSend	
k-1	recvReduceCopySend	
k to $2k-3$	recvCopySend	
2k-2	recv	

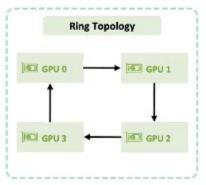


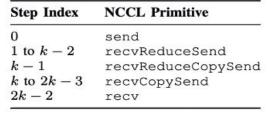


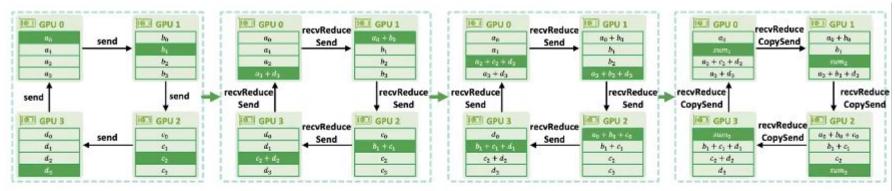
Step Index	NCCL Primitive	
0	send	
1 to $k-2$	recvReduceSend	
k-1	recvReduceCopySend	
k to 2k-3	recvCopySend	
2k-2	recv	

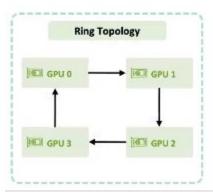
Step Index	NCCL Primitive	
0	send	
1 to $k-2$	recvReduceSend	
k-1	recvReduceCopySend	
k to 2k-3	recvCopySend	
2k-2	recv	



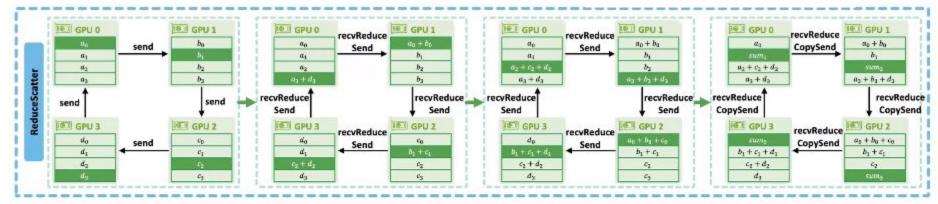


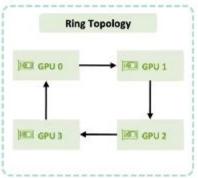






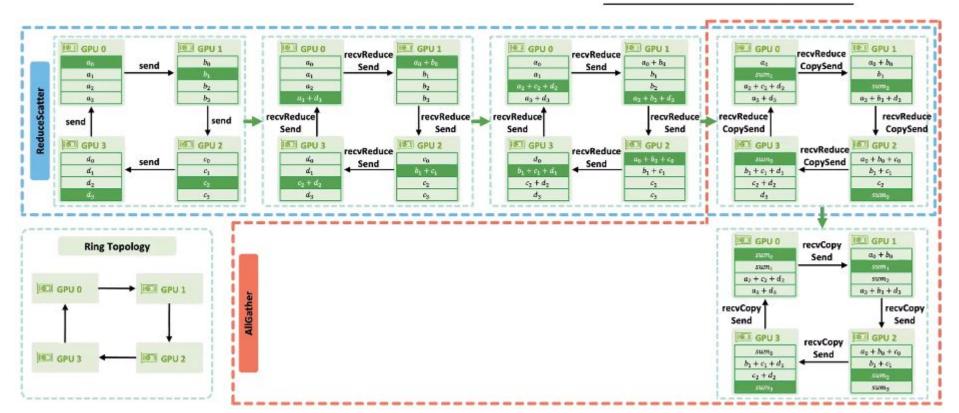
Step Index	NCCL Primitive	
0	send	
1 to $k-2$	recvReduceSend	
k-1	recvReduceCopySend	
k to $2k-3$	recvCopySend	
2k-2	recv	





Ring AllGather Example

Step Index	NCCL Primitive	
)	send	
1 to $k-2$	recvReduceSend	
k-1	recvReduceCopySend	
k to 2k-3	recvCopySend	
2k-2	recv	



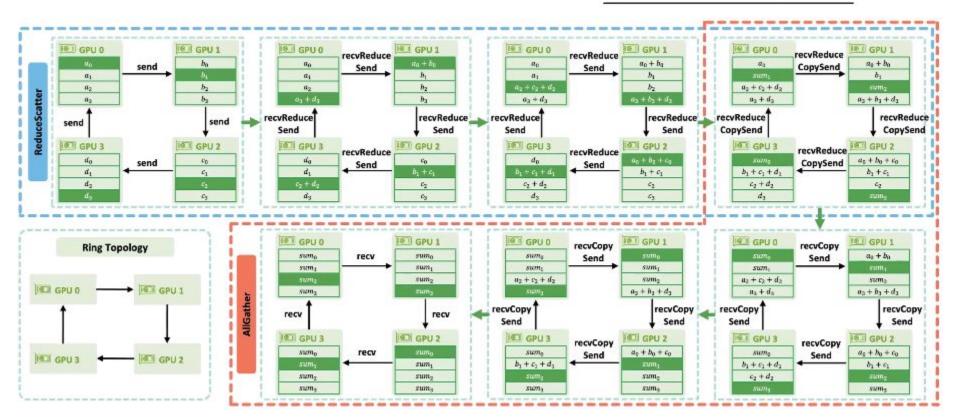
Ring AllGather Example

 $\begin{array}{lll} \textbf{Step Index} & \textbf{NCCL Primitive} \\ \hline 0 & & \texttt{send} \\ 1 \texttt{ to } k-2 & \texttt{recvReduceSend} \\ k-1 & & \texttt{recvReduceCopySend} \\ k \texttt{ to } 2k-3 & & \texttt{recvCopySend} \\ 2k-2 & & \texttt{recv} \end{array}$

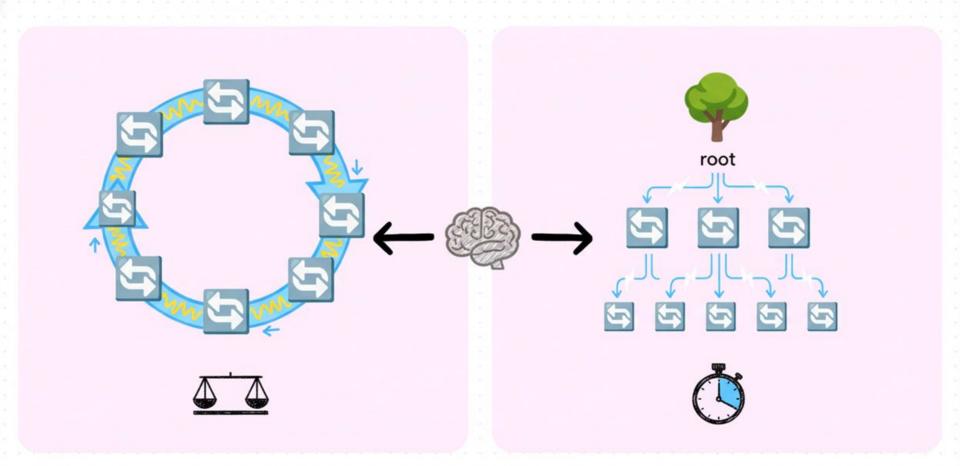


Ring AllGather Example

Step Index	NCCL Primitive	
0	send	
1 to $k-2$	recvReduceSend	
k-1	recvReduceCopySend	
k to 2k-3	recvCopySend	
2k-2	recv	



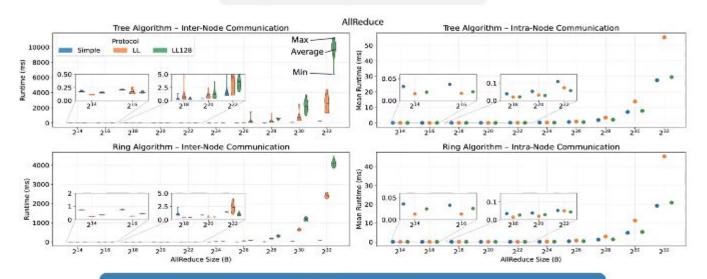
Choosing Between the Topologies



Benchmarking Results

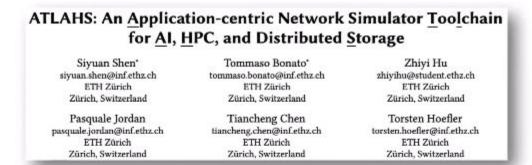
Alps Supercomputer (CSCS)

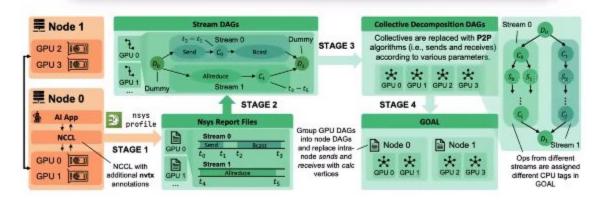
- Grace Hopper Superchips (GH200)
- 150 GB/s intra-node communication
- 25 GB/s Cray Slingshot
- · Dragonfly topology



Benchmarking results for all NCCL collectives in the paper

Impact and Outlook





Questions for Discussion

- How does NCCL manage communication channels, and what is the trade-off in channel count selection?
- NCCL uses different topologies (ring vs. tree) for different operations. How might emerging network architectures (like disaggregated systems, optical interconnects, or SmartNICs) fundamentally change optimal collective communication strategies?
- NCCL's "autotuning" approach generally provides good performance. How should we think about the balance between automatic optimization and manual control in HPC/ML systems?
- The benchmarking shows significant performance differences between intra-node and inter-node communication. What does this tell us about the future of large-scale AI training systems?
- How does the derived tool, ATLAHS, address the gap of NCCL being a black box in a way that simply accessing the source code cannot?