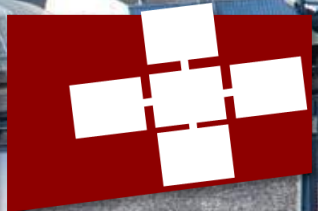


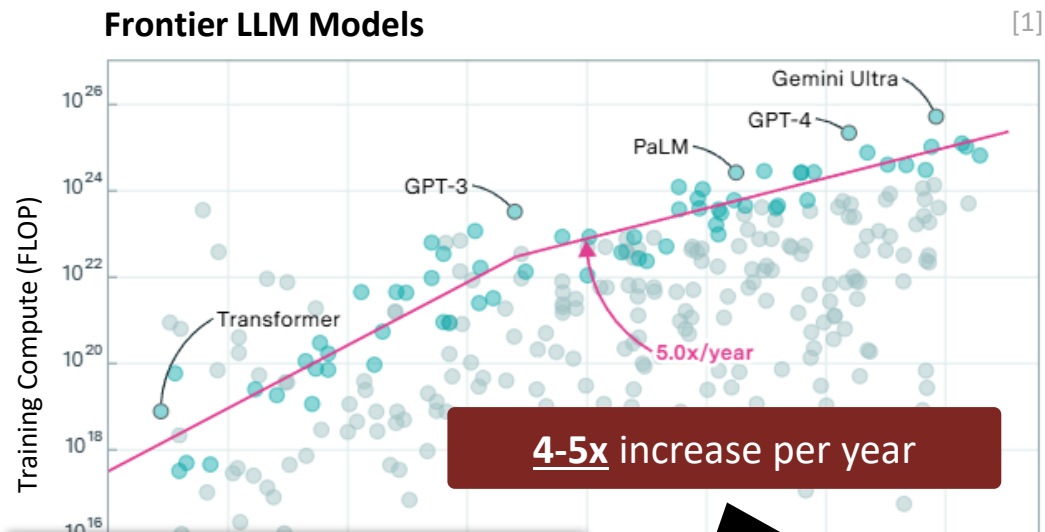
ATC'25, Boston, MA, USA

CrossPipe: Towards Optimal Pipeline Schedules for Cross-Datacenter Training

TIANCHENG CHEN, ALES KUBICEK, LANGWEN HUANG, TORSTEN HOEFLER



Motivation



The New York Times
Hungry for Energy, Amazon, Google and Microsoft Turn to Nuclear Power
 Large technology companies are investing in nuclear energy as an emissions-free source of power for their artificial intelligence and other business operations.
 IN-BRIEF ANALYSIS
 OCTOBER 1, 2024

Data center owners turn to nuclear as potential electricity source
 Nuclear power plants that have signed agreements to power data centers (as of Sep 2024)

Increasing power demand

HOME > NEWS > IT HARDWARE & SEMICONDUCTORS

Training Google's Gemini Ultra data centers, and rising power demand

What we can learn from research published by Google Cloud's CEO

September 4, 2024

Multi-Datacenter Training: OpenAI's Ambitious Plan To Beat Google's Infrastructure

// Gigawatt Clusters, Telecom Networking, Long Haul Fiber, Hierarchical & Asynchronous SGD, Distributed Infrastructure Winners

By Dylan Patel, Daniel Nishball and Jeremie Eliahou Ontiveros

34 minutes, 18 comments

"We don't disclose exactly the details of how many locations but Gemini Ultra was trained across multiple sites, and multiple clusters within those sites."

Thomas Kurian, CEO Google Cloud
 December 2023

Problem: local power limitation

Solution: training across multiple datacenters

[1] Training Compute of Frontier AI Models Grows by 4-5x per Year, Report, [Epoch AI, 2024](#)

Motivation




HOME > NEWS > IT HARDWARE & SEMICONDUCTORS

Training Google's Gemini data centers, and rising power demand

September 4, 2024

Multi-Datcenter Training: OpenAI's Ambitious Plan To Beat Google's Infrastructure

// Gigawatt Clusters, Telecom Networking, Long Haul Fiber, Hierarchical & ...



How do we do cross-datacenter training efficiently with open-source stack?

Increasing power demand

Hungry for Energy, Amazon, Google and Microsoft Turn to Nuclear Power

Large technology companies are investing in nuclear energy as an emissions-free power source for artificial intelligence and other business operations.

OCTOBER 1, 2024

Data center owners turn to nuclear as potential electricity source

Nuclear power plants that have signed agreements to power data centers (as of Sep 2024)

When you purchase through links on our site, we may earn a commission.

By Dylan Patel, Daniel Nightball and Jerome Eliahou Ontiveros

"We don't disclose exactly the details of how many locations but Gemini Ultra was trained across multiple sites, and multiple clusters within those sites."

Thomas Kurian, CEO Google Cloud
December 2023

Problem: local power limitation

Solution: training across multiple datacenters

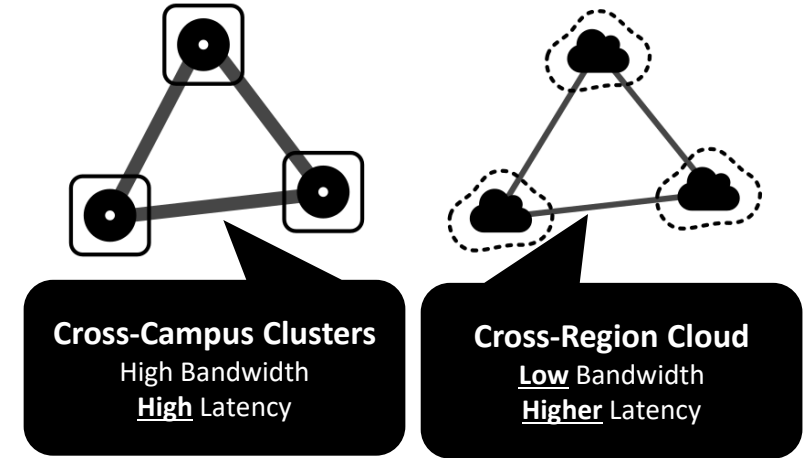
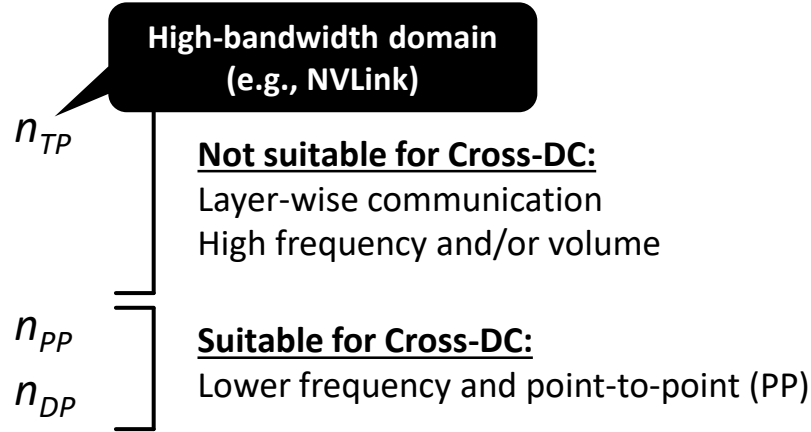
[1] Training Compute of Frontier AI Models Grows by 4-5x per Year, Report, [Epoch AI, 2024](#)

Problem and Parallelism Strategies

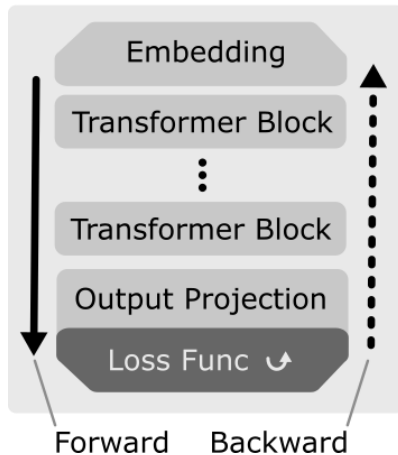
- **Challenge:** high latency / low bandwidth on cross-DC boundary

- Hybrid parallelism

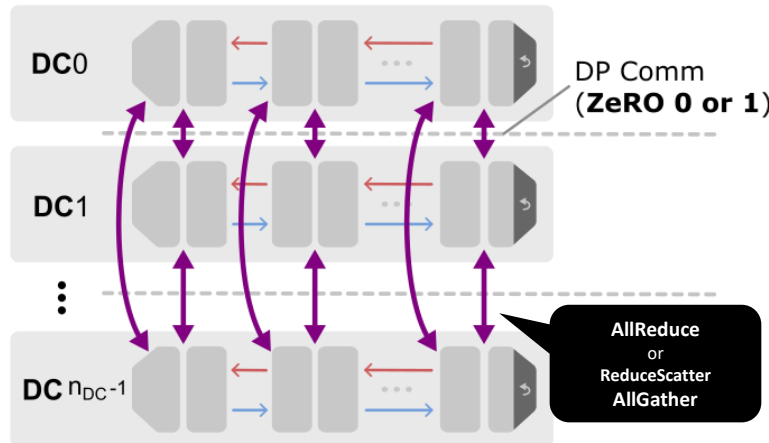
- Tensor (= Operator)
- Sequence (= Context)
- Expert
- **Pipeline**
- **Data**



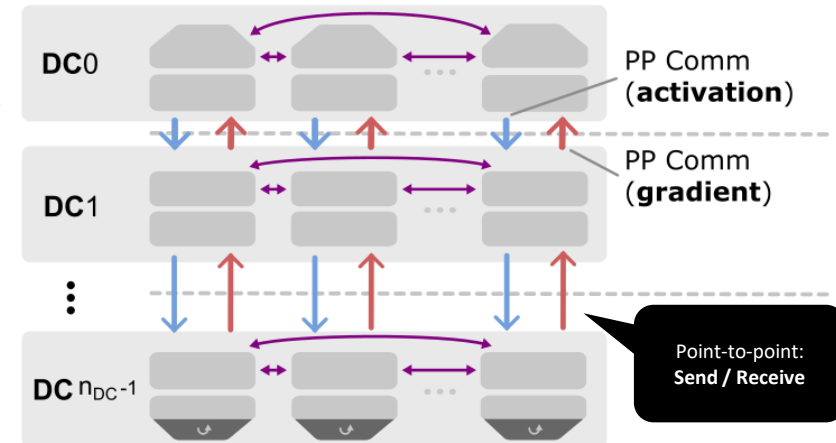
Model:



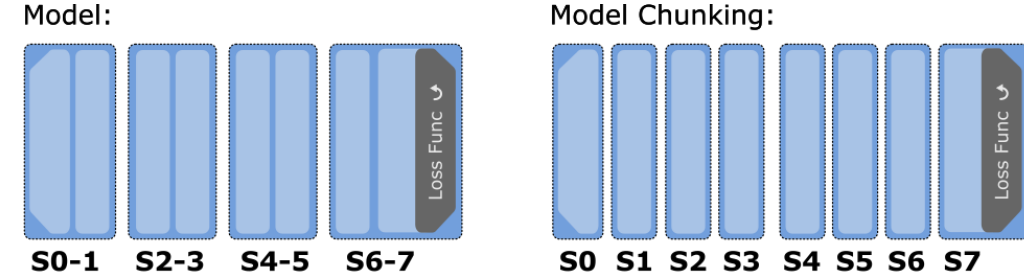
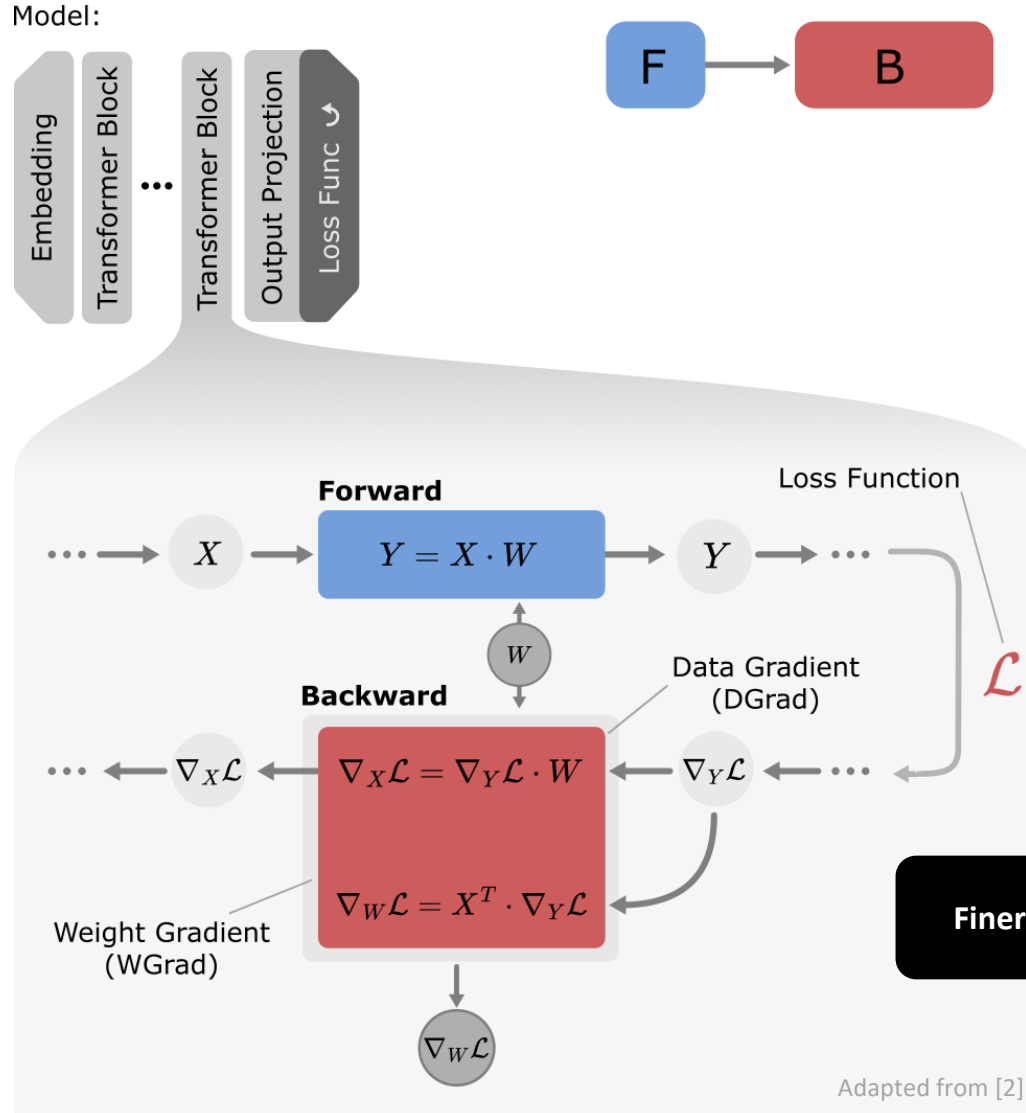
Cross-DC Data Parallelism (DP)



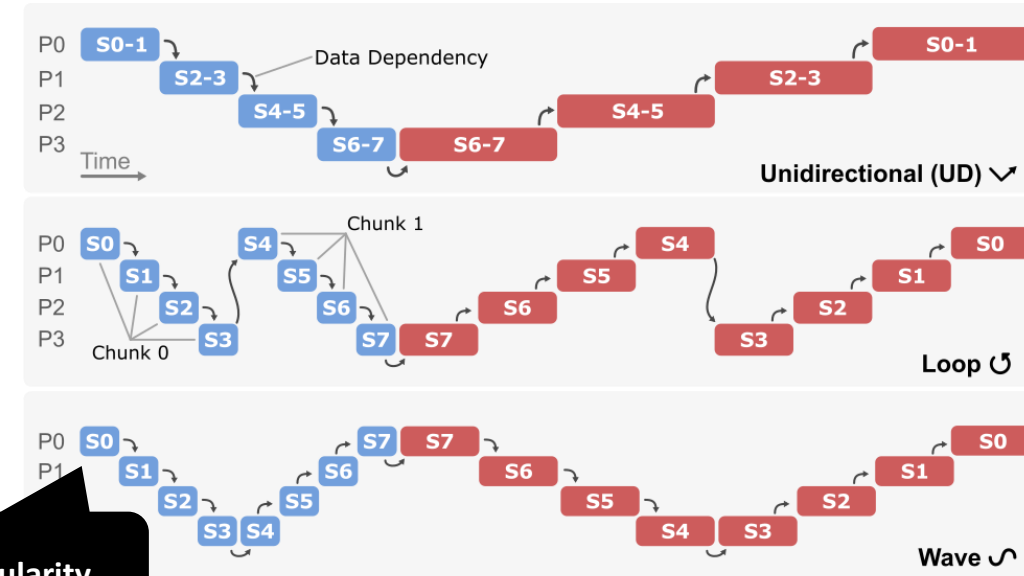
Cross-DC Pipeline Parallelism (PP)



Pipeline Parallelism – Data Dependencies

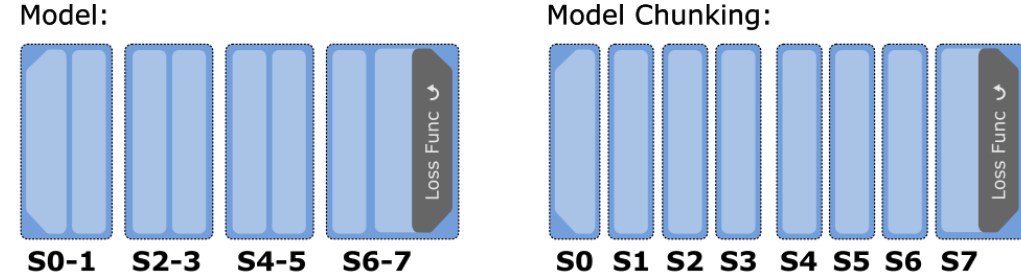
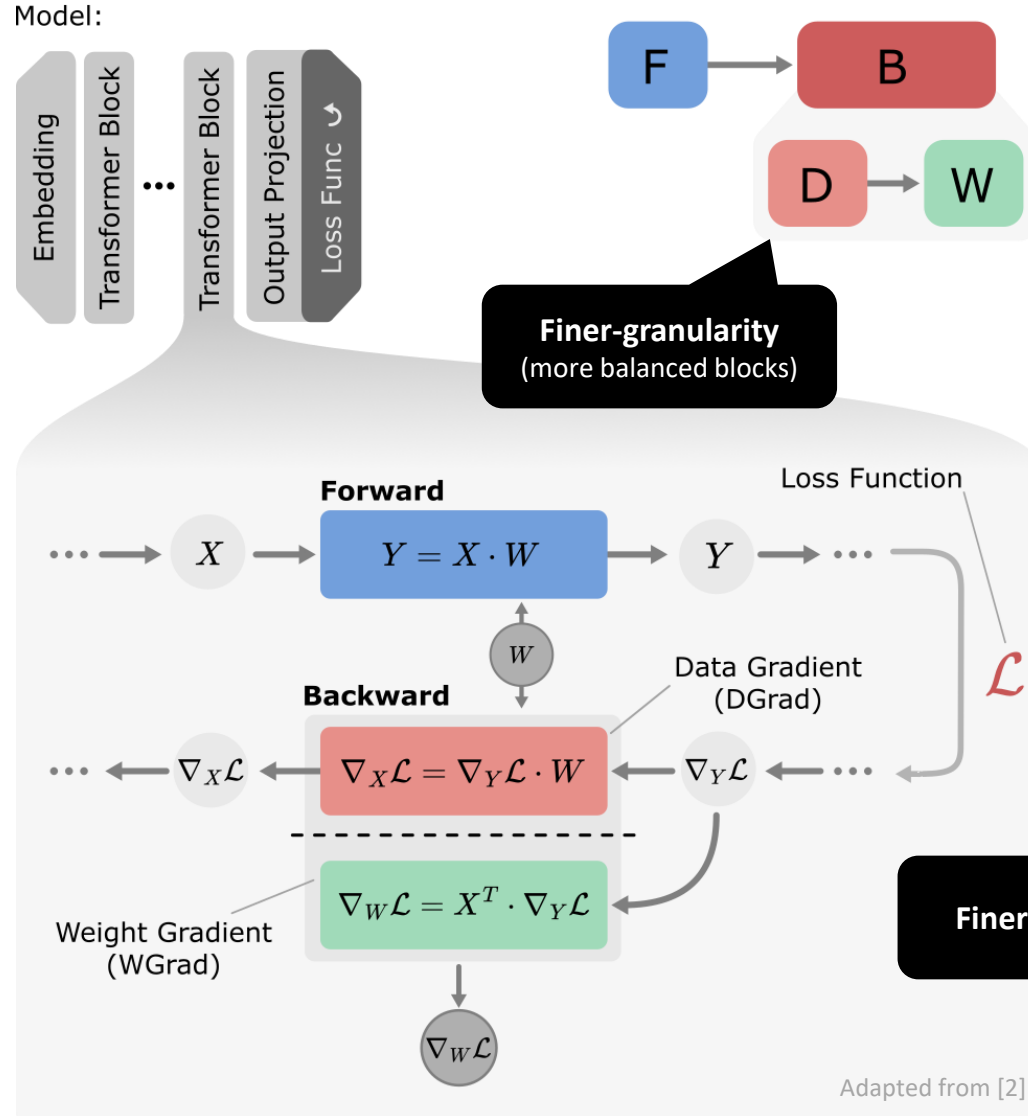


Traversal Patterns:

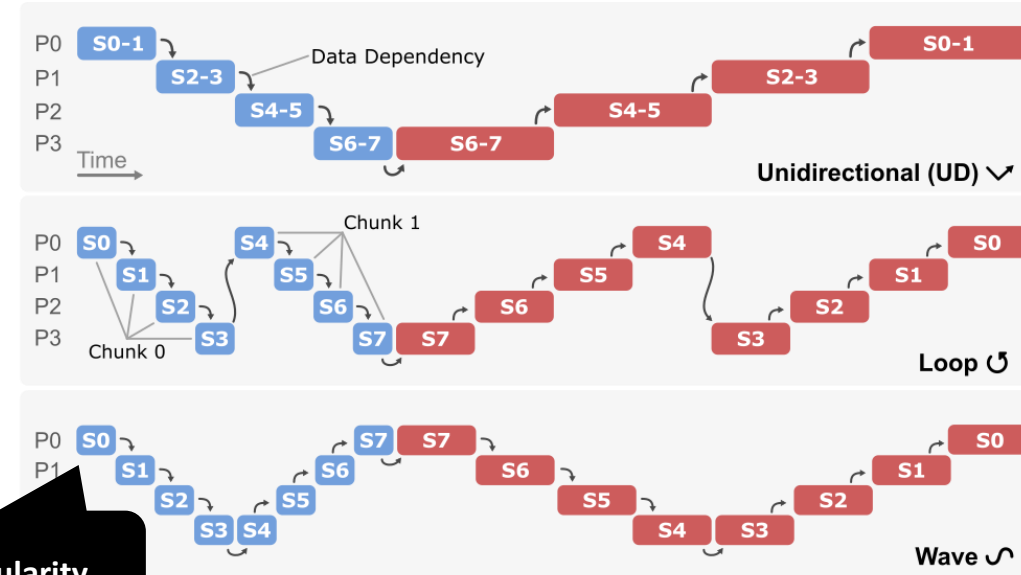


Finer-granularity

Pipeline Parallelism – Data Dependencies

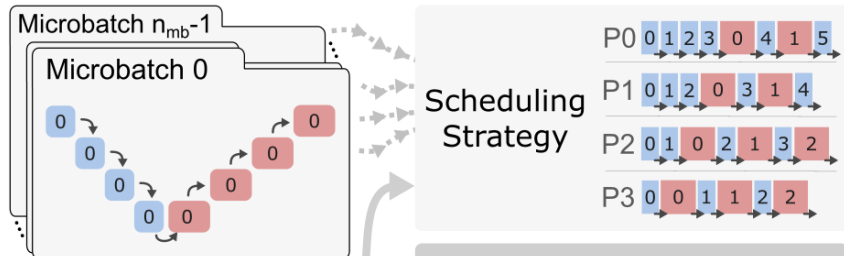


Traversal Patterns:



True dependencies at varying granularity

Pipeline Parallelism – Schedules



Data Dependencies

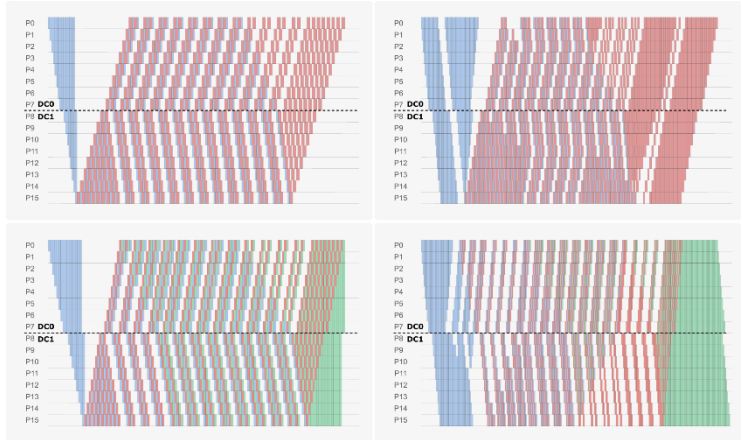
Traversal Patterns:

- Unidirectional (UD)
- Loop
- Wave

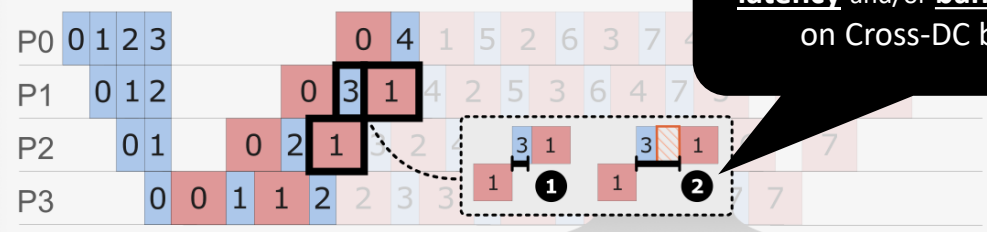
Schedule Dependencies

Pipeline Schedule

- 1F1B
- IV1F1B
- ZBH1
- ZBV

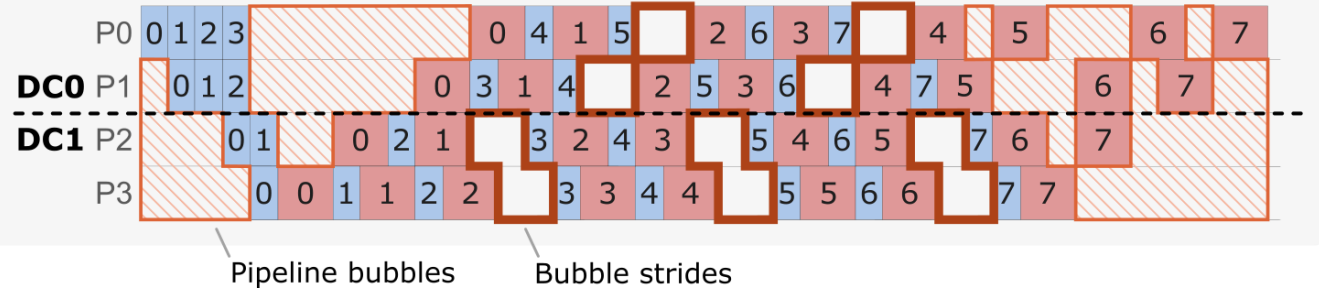


1 Single-DC - communication cost ignored



Significant latency and/or bandwidth delay on Cross-DC boundary

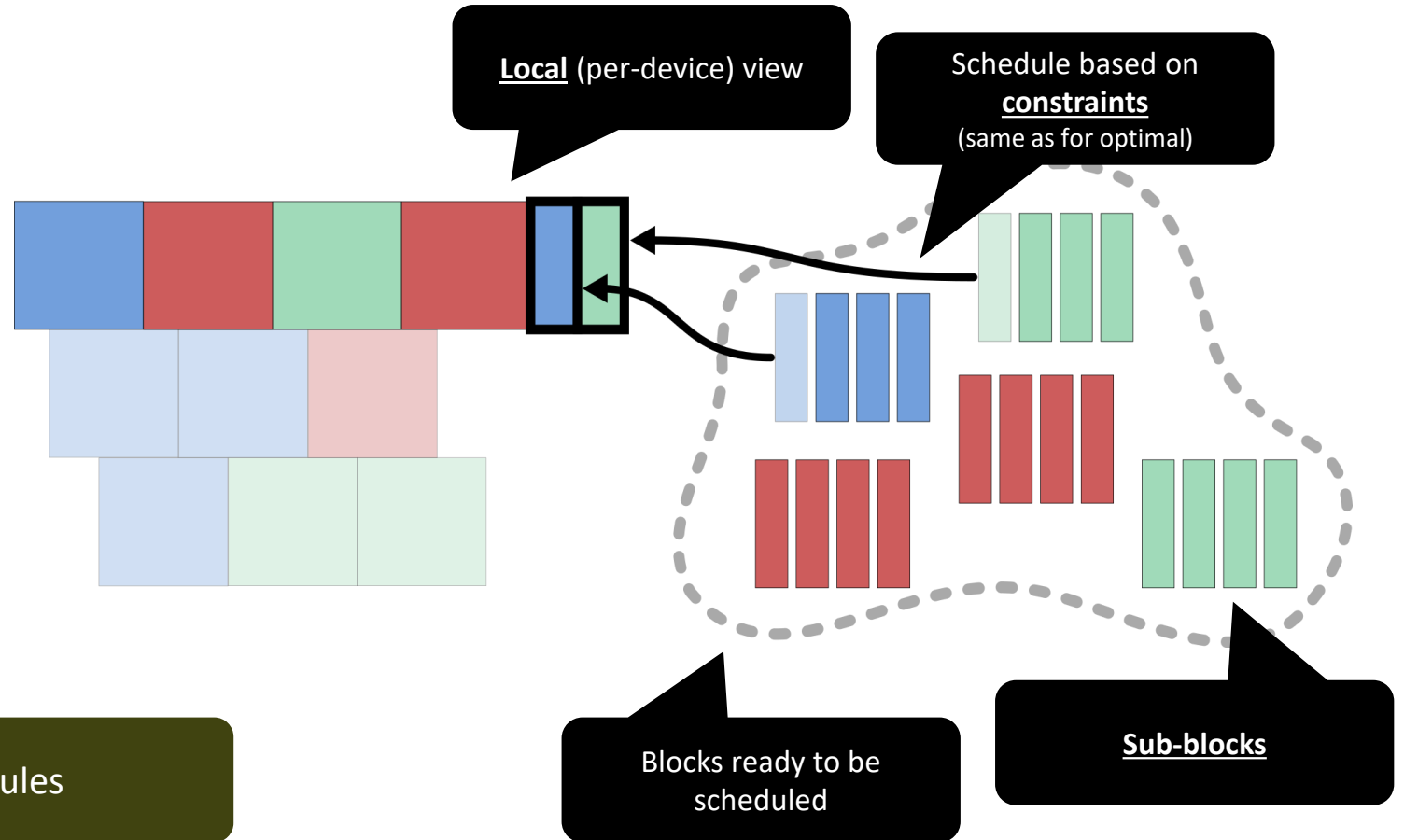
2 Cross-DC



What can we do about the bubble strides?

CrossPipe – Greedy Schedule

- Dynamic conditions (schedule adaptation)
- Solver runtime \Rightarrow limiting
- Algorithmic approach**
 \Rightarrow sub-block scheduling
- Faster schedule generation
- Same constraints
- Near-optimal



Attractive alternative to optimal schedules

Cross-DC DP vs. Cross-DC PP

Performance model

Model: Llama 3 405B

Based on Llama 3 technical report

- 2 DC
- $n_{TP} = 8$ $n_{PP} = 16$ $n_{DP} = 64$
- $T_F = \underline{109 \text{ ms}}$

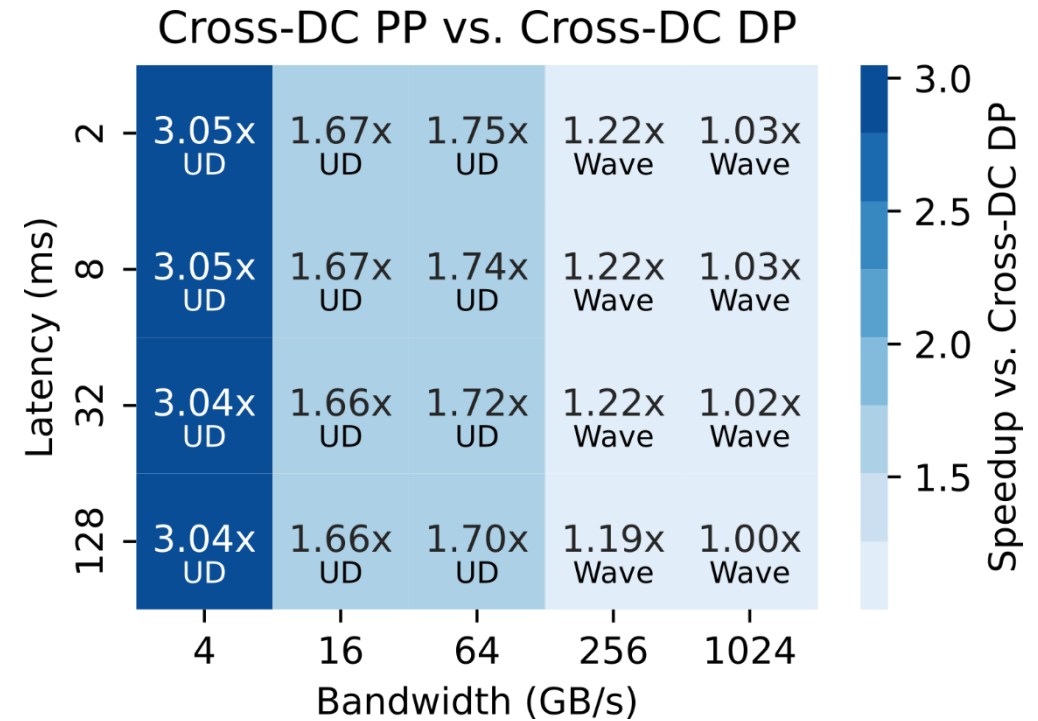
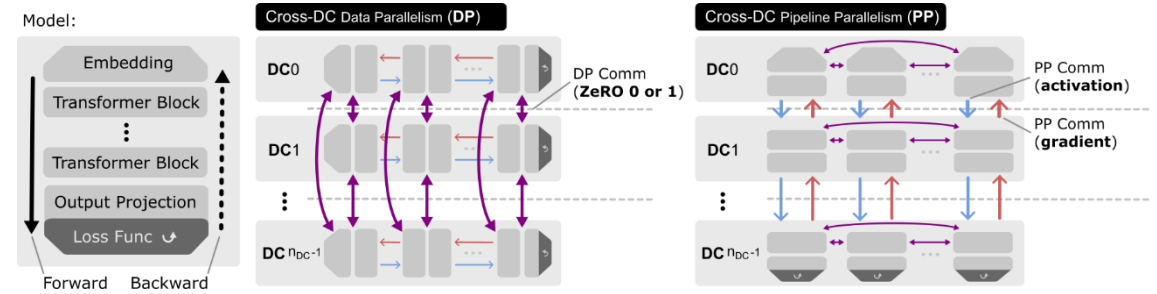
Time of forward (F) block

PP schedules:

Traversal pattern (best)

- Cross-DP: ZBV
- Cross-PP: CrossPipe (UD or Wave)

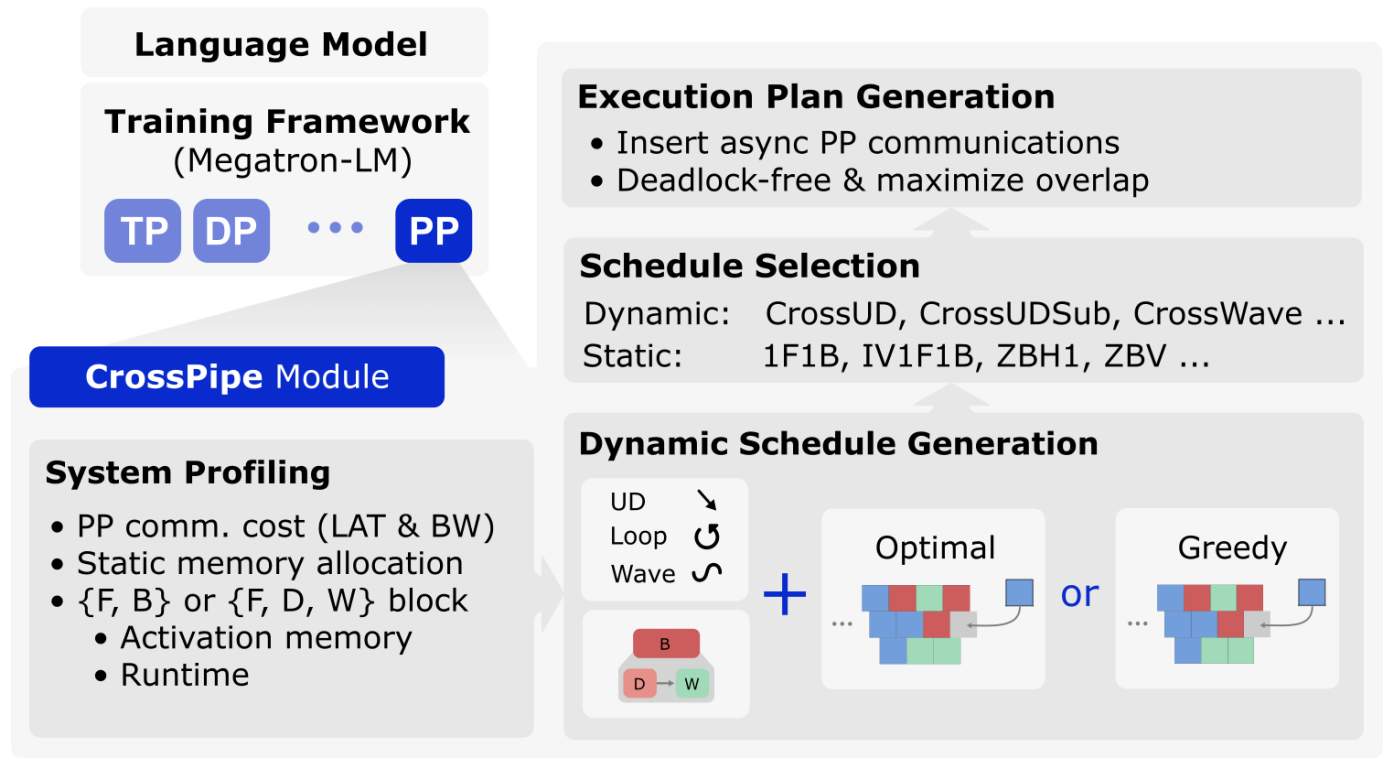
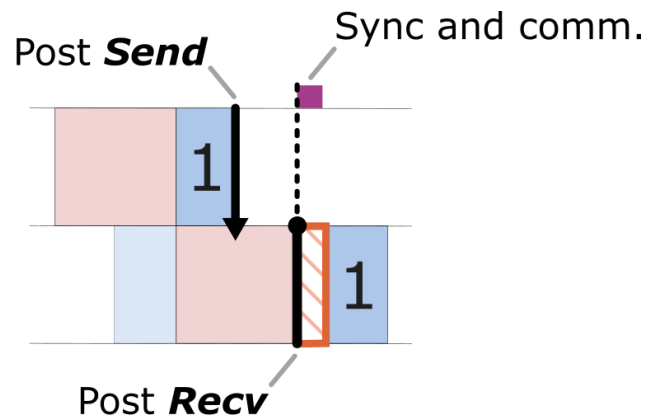
Cross-DC – Pipeline Parallelism → better choice



CrossPipe – Implementation

- **Optimization**
- Post *Recv* ahead of *Send*
 ⇒ Based on schedule – maximize overlap

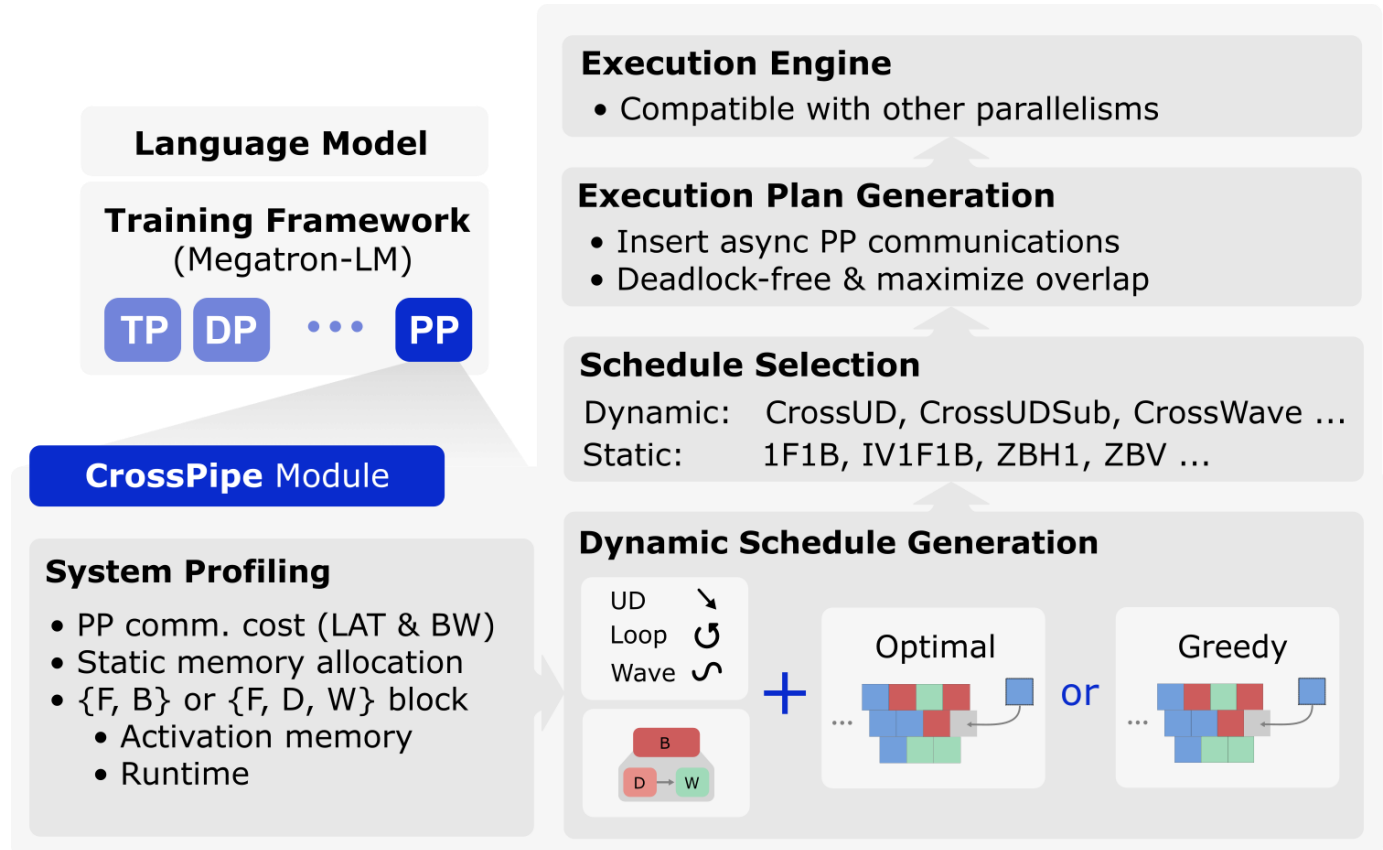
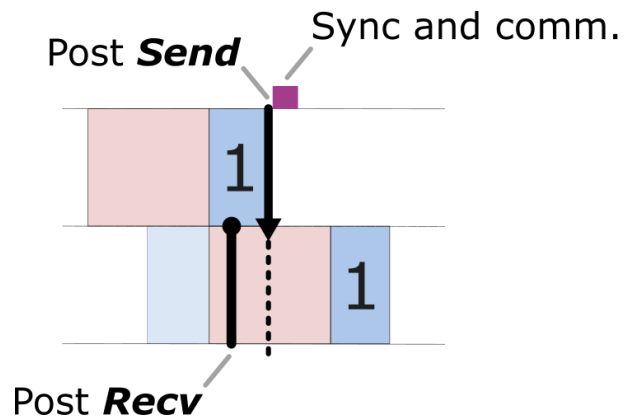
Original (NCCL):



CrossPipe – Implementation


- **Optimization**
- Post *Recv* ahead of *Send*
⇒ Based on schedule – maximize overlap
- Both **static** and **dynamic** schedules


Improved (NCCL):




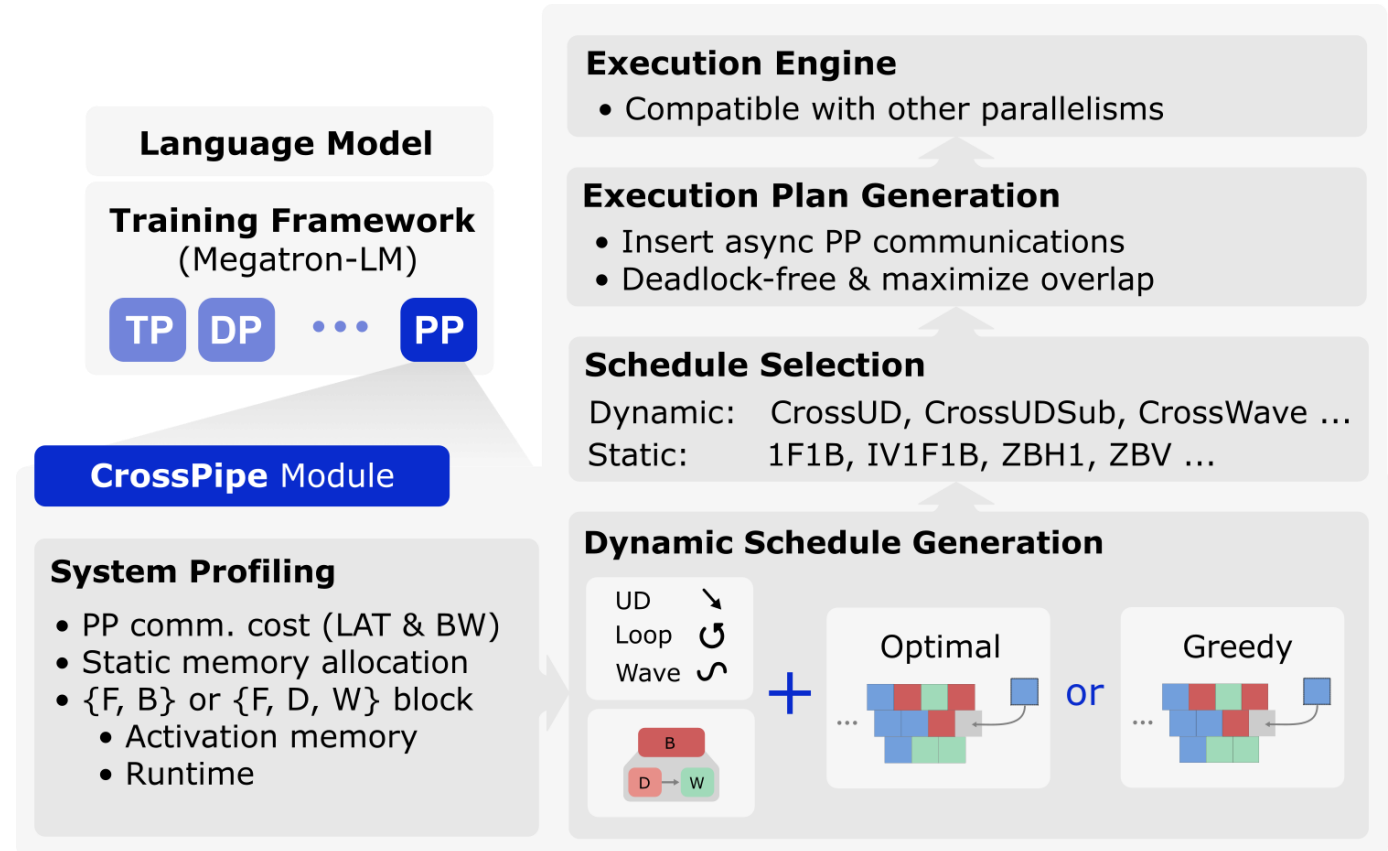
CrossPipe – Implementation

- **Optimization**
- Post *Recv* ahead of *Send*
⇒ Based on schedule – maximize overlap
- Both **static** and **dynamic** schedules


Adaptable (Single-DC / Cross-DC)


Extensible (configure / abstraction)


Efficient (optimal / near-optimal)

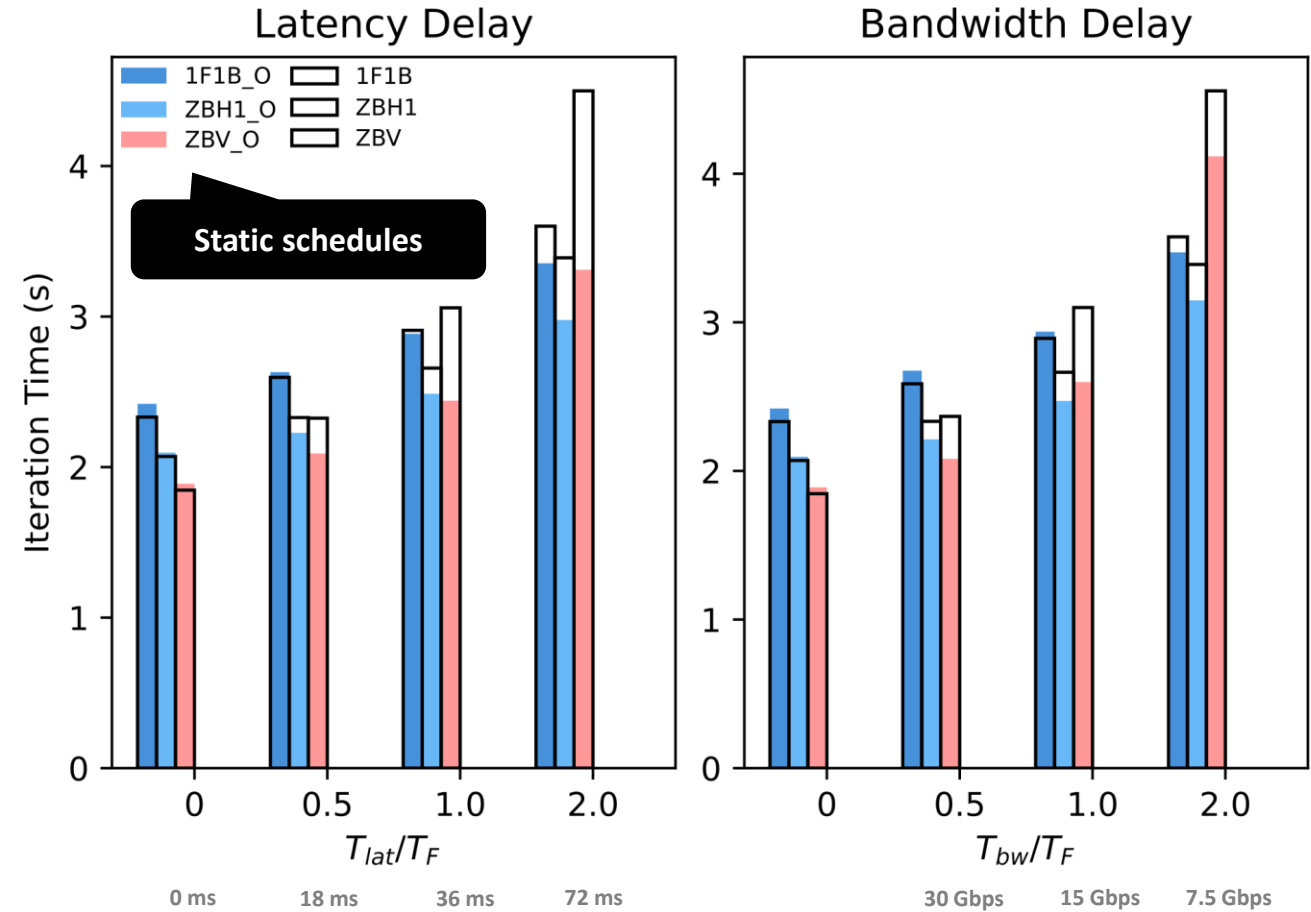


Evaluation

- **Single-DC system: Alps**
 - Node: 4x GH200 (96 GB HBM3), NVLink 4.0
 - HPE Cray Slingshot-11 interconnect
HPE Cray Cassini-1 200 Gb/s NIC | Dragonfly

- **Cross-DC latency and bandwidth delay**
 - Emulated in single-DC environment
 - PyTorch – injected spin kernels (validated)

- **Model: M70** (Llama 70B)
 - 2 DC
 - $n_{TP} = 4$ $n_{PP} = 8$ $n_{DP} = 1$
 - $T_F = \underline{36 \text{ ms}}$

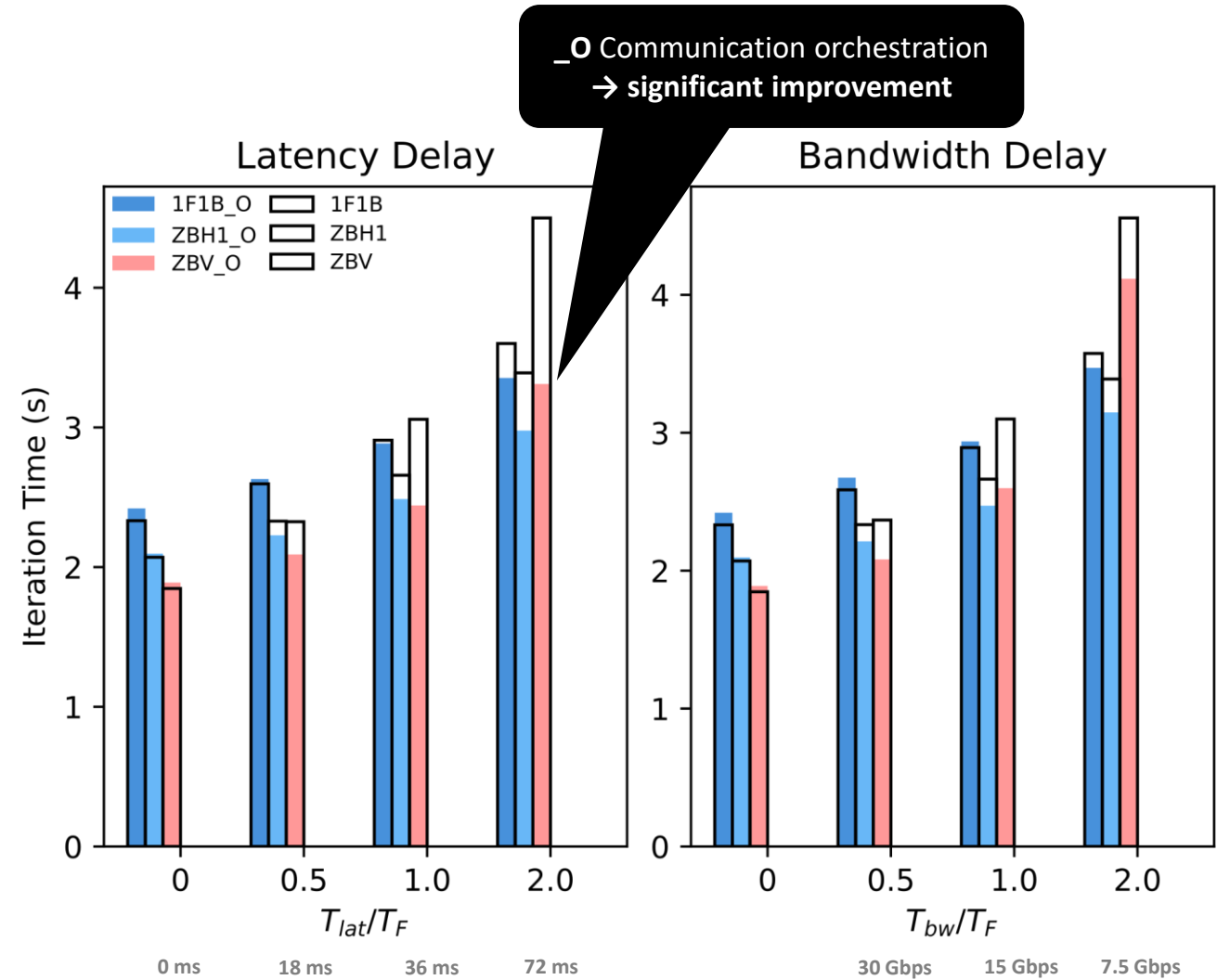


Evaluation

- **Single-DC system: Alps**
 - Node: 4x GH200 (96 GB HBM3), NVLink 4.0
 - HPE Cray Slingshot-11 interconnect
HPE Cray Cassini-1 200 Gb/s NIC | Dragonfly

- **Cross-DC latency and bandwidth delay**
 - Emulated in single-DC environment
 - PyTorch – injected spin kernels (validated)

- **Model: M70** (Llama 70B)
 - 2 DC
 - $n_{TP} = 4$ $n_{PP} = 8$ $n_{DP} = 1$
 - $T_F = \underline{36\text{ ms}}$

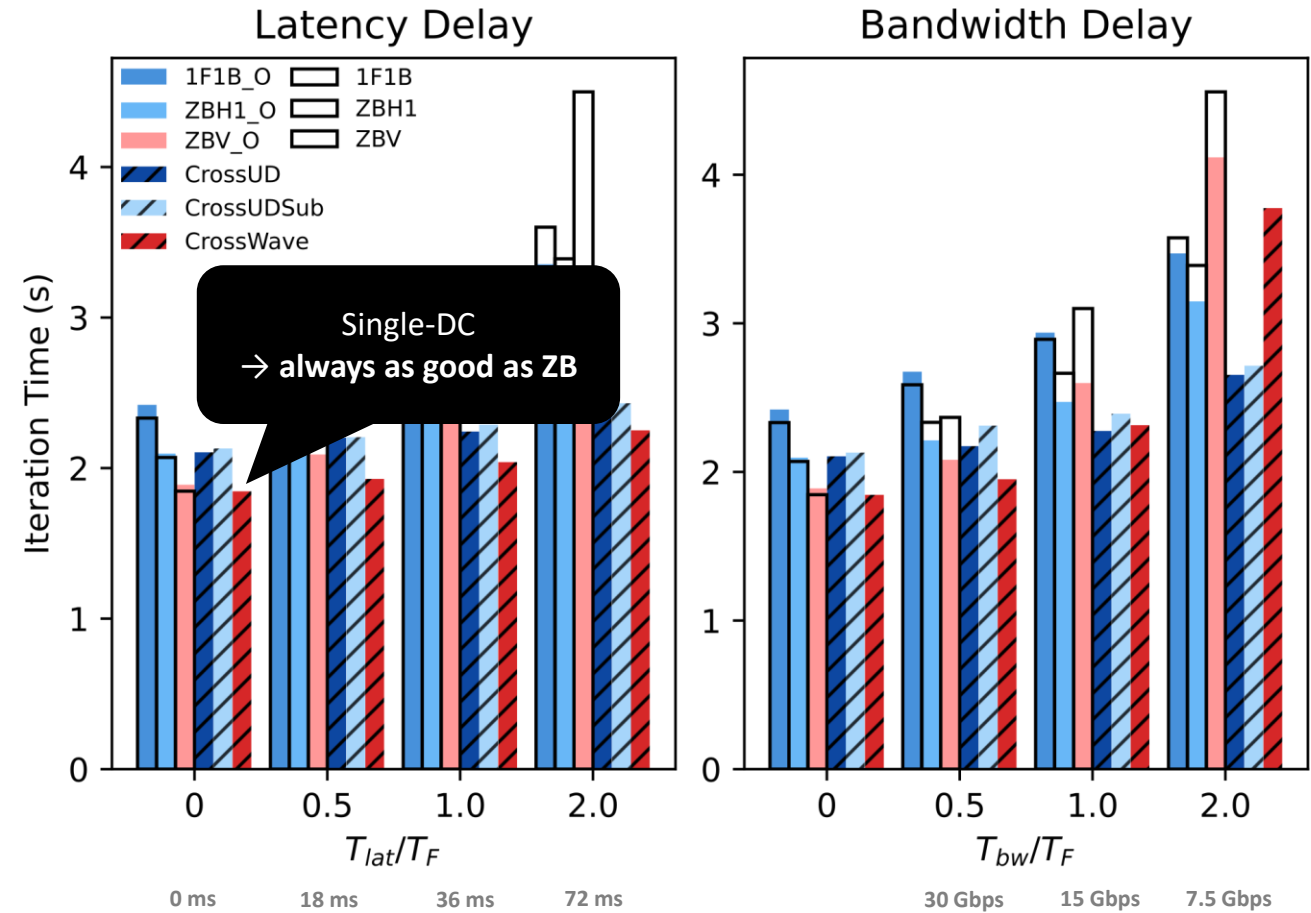


Evaluation

- **Single-DC system: Alps**
 - Node: 4x GH200 (96 GB HBM3), NVLink 4.0
 - HPE Cray Slingshot-11 interconnect
HPE Cray Cassini-1 200 Gb/s NIC | Dragonfly

- **Cross-DC latency and bandwidth delay**
 - Emulated in single-DC environment
 - PyTorch – injected spin kernels (validated)

- **Model: M70 (Llama 70B)**
 - 2 DC
 - $n_{TP} = 4$ $n_{PP} = 8$ $n_{DP} = 1$
 - $T_F = \underline{36\text{ ms}}$

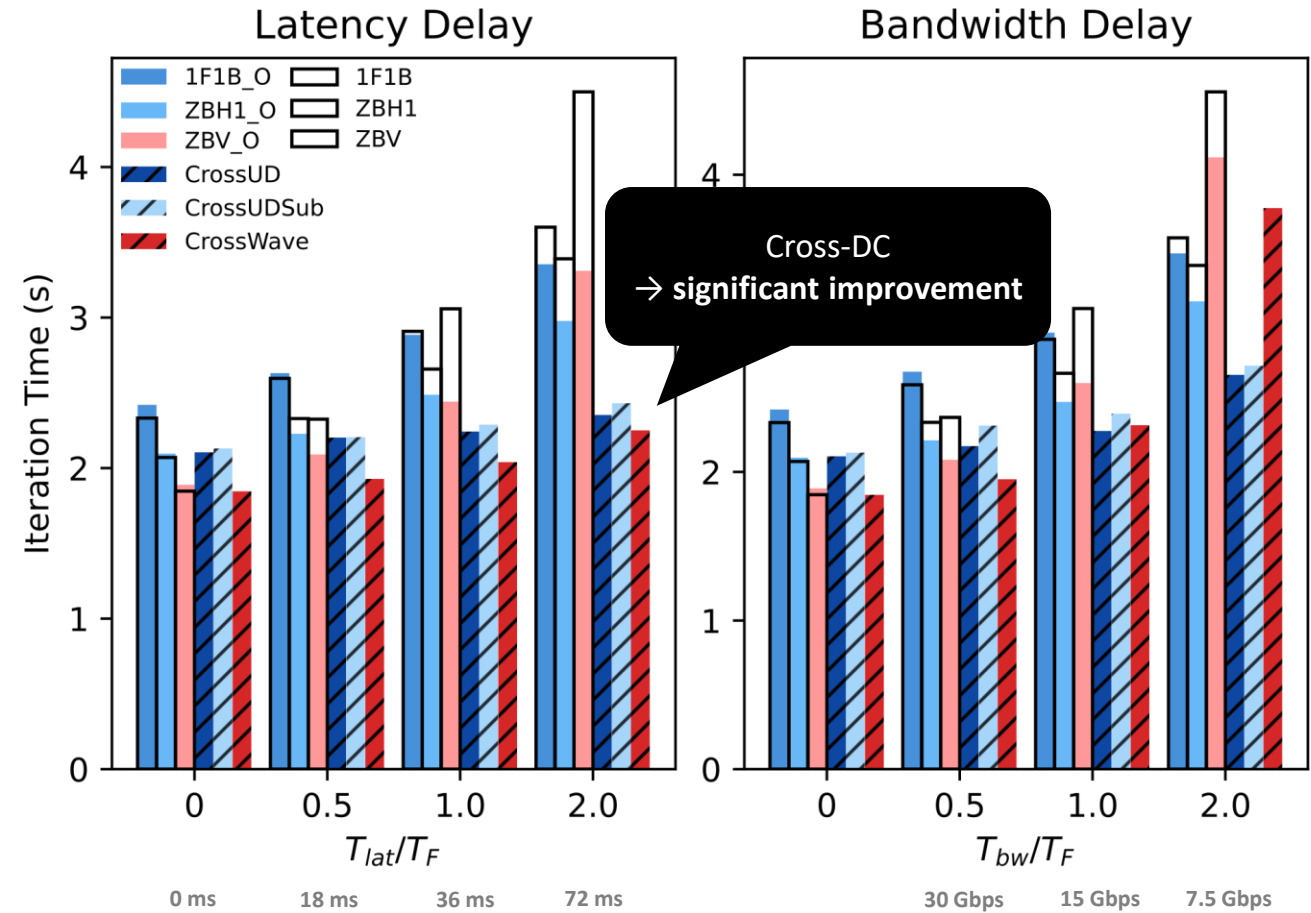


Evaluation

- **Single-DC system: Alps**
 - Node: 4x GH200 (96 GB HBM3), NVLink 4.0
 - HPE Cray Slingshot-11 interconnect
HPE Cray Cassini-1 200 Gb/s NIC | Dragonfly

- **Cross-DC latency and bandwidth delay**
 - Emulated in single-DC environment
 - PyTorch – injected spin kernels (validated)

- **Model: M70** (Llama 70B)
 - 2 DC
 - $n_{TP} = 4$ $n_{PP} = 8$ $n_{DP} = 1$
 - $T_F = \underline{36\text{ ms}}$

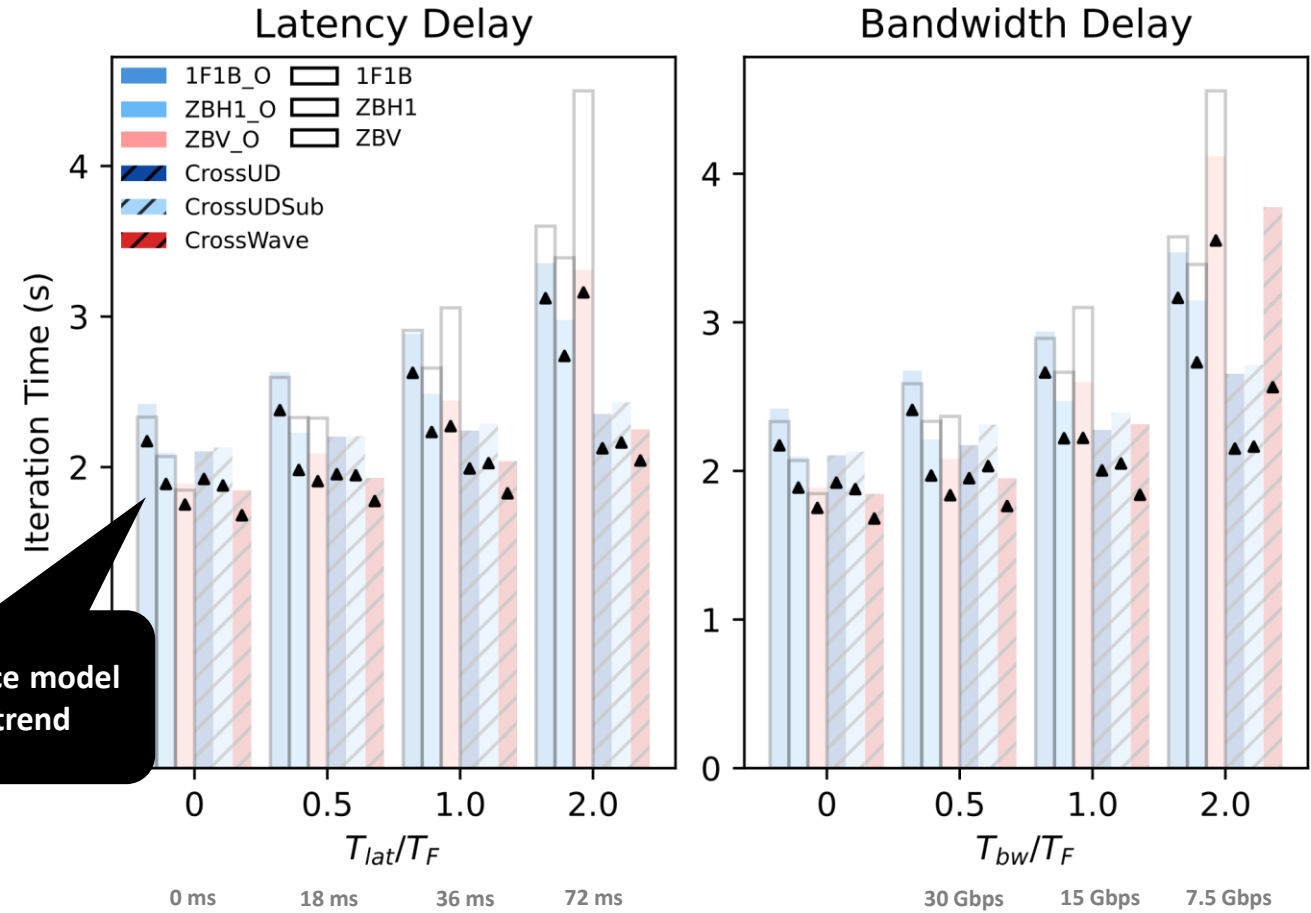


Evaluation

- **Single-DC system: Alps**
 - Node: 4x GH200 (96 GB HBM3), NVLink 4.0
 - HPE Cray Slingshot-11 interconnect
HPE Cray Cassini-1 200 Gb/s NIC | Dragonfly

- **Cross-DC latency and bandwidth delay**
 - Emulated in single-DC environment
 - PyTorch – injected spin kernels (validated)

- **Model: M70 (Llama 70B)**
 - 2 DC
 - $n_{TP} = 4$ $n_{PP} = 8$ $n_{DP} = 1$
 - $T_F = \underline{36\text{ ms}}$

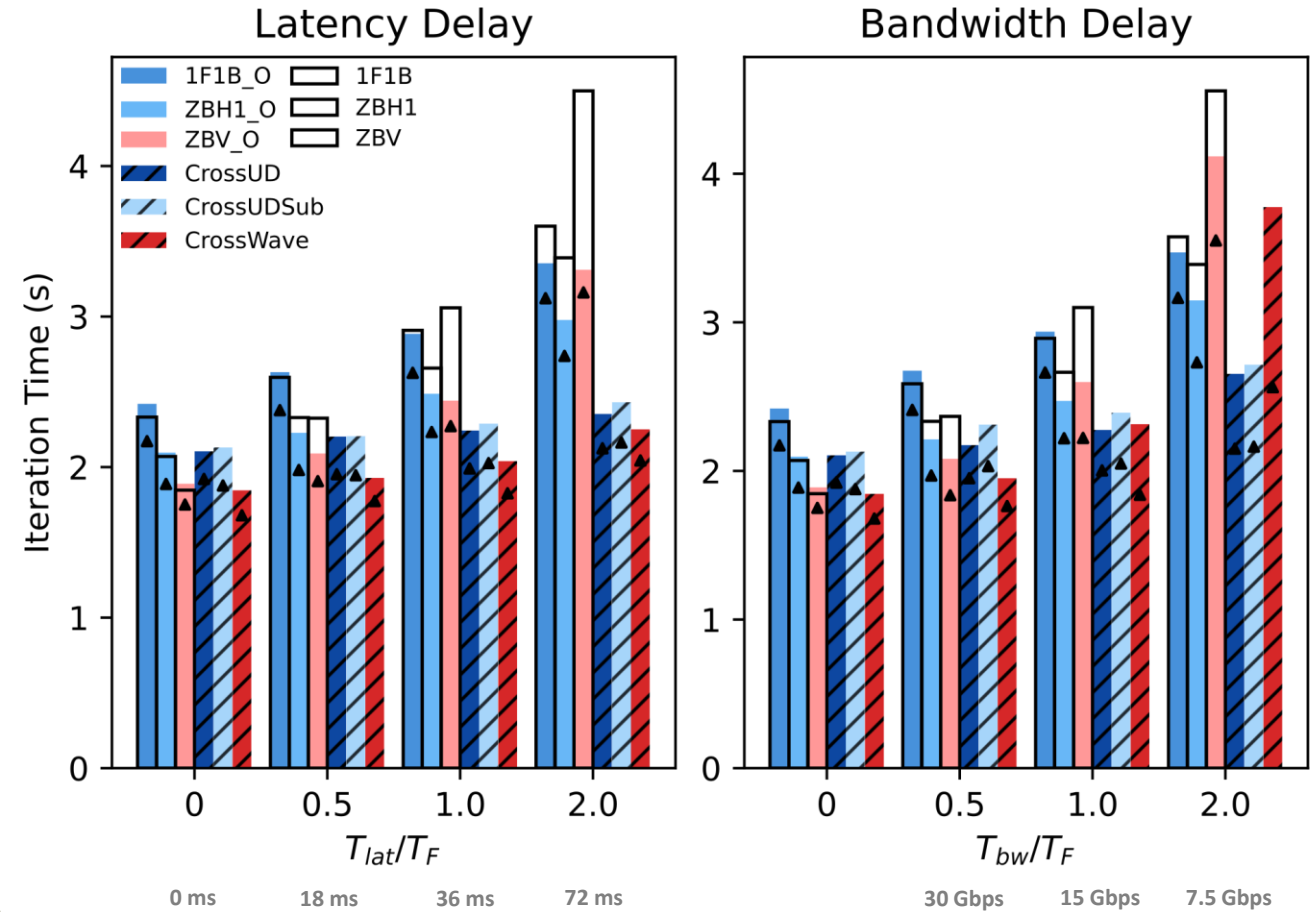


Evaluation

- **Single-DC system: Alps**
 - Node: 4x GH200 (96 GB HBM3), NVLink 4.0
 - HPE Cray Slingshot-11 interconnect
HPE Cray Cassini-1 200 Gb/s NIC | Dragonfly

- **Cross-DC latency and bandwidth delay**
 - **Emulated** in single-DC environment
 - PyTorch – injected spin kernels (validated)

- **Model: M70** (Llama 70B)
 - 2 DC
 - $n_{TP} = 4$ $n_{PP} = 8$ $n_{DP} = 1$
 - $T_F = \underline{36\text{ ms}}$



Up to **33.6%** training time reduction
(identical memory constraints)

Evaluation

Multiple Models: M8 and M70

Varying number of DCs: 2 and 4

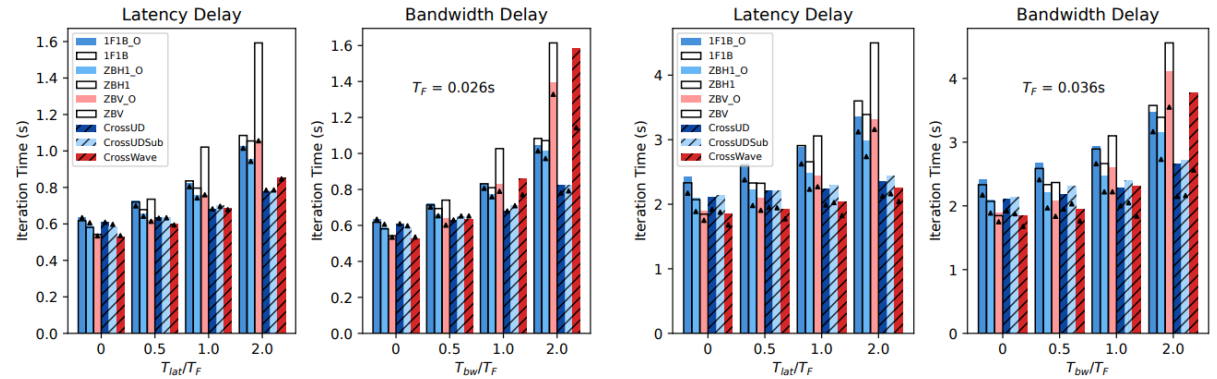
Increased GBS

Increased memory budget

Activation recompute

PP and DP tradeoff

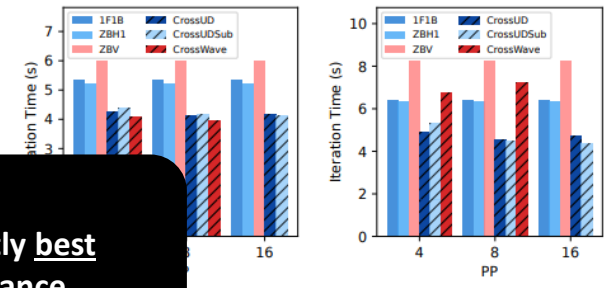
Mix latency + bandwidth delay



(a) Model M8

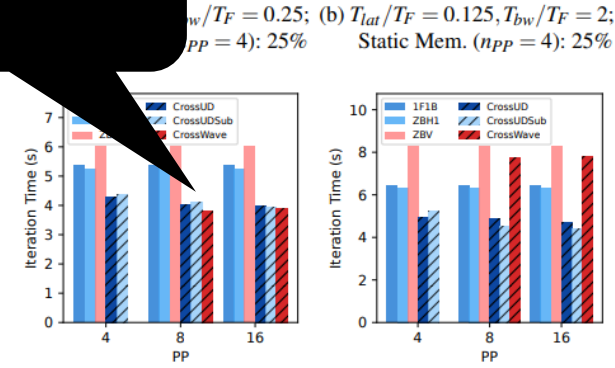
(b) Model M70

T_{lat}/T_F	T_{bw}/T_F	Case	Static			Dynamic (This Work)		
			1F1B	ZBH1	ZBV	UDSub	UD	Wave
0	0	1	0.151	0.133	0.118	0.137	0.137	0.119
		2	0.174	0.168	0.161	-	0.165	0.157
		3	-	-	-	-	0.121	0.118
0.25	0.25	1	0.168	0.15	0.148	0.149	0.142	0.127
		2	0.193	0.187	0.177	-	0.17	0.159
		3	-	-	-	0.188	0.188	0.118
0.25	2	1	0.241	0.23	0.315	-	-	-
		2	0.262	0.259	0.33	-	-	-
		3	-	-	-	-	-	-
2	0.25	1	0.242	0.229	0.314	-	-	-
		2	0.262	0.258	0.329	-	-	-
		3	-	-	-	-	-	-
2	2	1	0.321	0.309	0.473	-	-	-
		2	0.333	0.331	0.476	-	-	-
		3	-	-	-	-	-	-



Consistently best performance

T_{lat}/T_F	T_{bw}/T_F	Case	Static			Dynamic (This Work)		
			1F1B	ZBH1	ZBV	UDSub	UD	Wave
0	0	1	0.149	0.133	0.119	0.138	0.138	0.122
		2	0.173	0.168	0.16	-	0.167	0.157
		3	-	-	-	0.123	0.119	0.115
0.25	0.25	1	0.177	0.158	0.161	0.155	0.148	0.141
		2	0.198	0.19	0.181	-	0.173	0.163
		3	-	-	-	0.126	0.123	0.115
0.25	2	1	0.269	0.249	0.339	0.216	0.217	0.286
		2	0.274	0.269	0.331	-	0.198	0.29
		3	-	-	-	-	0.162	0.262
2	0.25	1	0.268	0.248	0.337	0.2	0.197	0.214
		2	0.274	0.269	0.33	-	0.184	0.177
		3	-	-	-	0.138	0.138	0.139
2	2	1	0.359	0.338	0.512	0.268	0.271	0.339
		2	0.349	0.346	0.479	-	0.213	0.295
		3	-	-	-	0.178	0.178	0.264



(c) $T_{lat}/T_F = 1, T_{bw}/T_F = 0.25$; (d) $T_{lat}/T_F = 0.125, T_{bw}/T_F = 2$; Static Mem. ($n_{PP} = 4$): 50%

Conclusions

Cross-DC DP vs. Cross-DC PP

- Performance model**
 - Model: Llama 3 405B
 - 2 DC
 - $\rho_{DP} = 8$ $\rho_{PP} = 16$ $\rho_{DP} = 64$
 - $T_f = 109$ ms
- PP schedules:**
 - Cross-DC: ZBV
 - Cross-PP: CrossPipe (UD or Wave)

Based on Llama 3 technical report

Traversal pattern (best)

Cross-DC - Pipeline Parallelism → better choice

Bandwidth (GB/s)	Latency (ms)	Speedup vs. Cross-DC DP
4	3.05x UD	~2.5
16	1.67x UD	~2.0
64	1.75x UD	~1.8
256	1.22x Wave	~1.5
1024	1.03x Wave	~1.4

CrossPipe – Optimal Schedule

CrossPipe – Greedy Schedule

Objective: minimize

Reorder

Smaller feasible ratio

Faster

Same peak memory

Optimal

Near-optimal

More of SPCL's research:

youtube.com/@spcl **210+ Talks**

twitter.com/spcl_eth **1.6K+ Followers**

github.com/spcl **5.6K+ Stars**

... or spcl.ethz.ch



CrossPipe – Implementation

- Communication orchestration**
 - NCCL – 4 GPU streams (Send, Recv) × (Next, Prev)
 - Post Recv ahead of Send → Based on profiling – maximize overlap
 - Both static and dynamic schedules
 - Improves static schedules in cross-DC

Adaptable (Single DC / Cross DC)

Extensible (configure / schedules)

Efficient (optimal)

Language Model

Training Framework (Megatron-LM)

Execution Engine

- Compatible with other parallelisms

Execution Plan Generation

- Insert async PP communications
- Deadlock-free & maximize overlap

Schedule Selection

Dynamic: CrossUD, CrossUDSub, CrossWave ...

Static: 1F1B, 1V1F1B, ZBH1, ZBV ...

Dynamic Schedule Generation

UD, Loop, Wave

Optimal or Greedy

System Profiling

- PP comm. cost (LAT & BW)
- Static memory allocation
- (F, B) or (F, D, W) block
- Activation memory
- Runtime

Evaluation

- Multiple Models: M8 and M70
- Varying number of DCs: 2 and 4
- Increased GBS
- Increased memory budget
- Activation recompute
- PP and DP tradeoff
- Mix latency + bandwidth delay

(a) Model M8

(b) Model M70

(c) $\rho_{DP} = 1, \rho_{PP} / \rho_{DP} = 0.25$ (d) $\rho_{DP} / \rho_{PP} = 0.125, \rho_{PP} / \rho_{DP} = 2$

(e) $\rho_{DP} = 1, \rho_{PP} / \rho_{DP} = 0.25$ (f) $\rho_{DP} / \rho_{PP} = 0.125, \rho_{PP} / \rho_{DP} = 2$

Open source:

github.com/spcl/crosspipe

