





Mist: Efficient Distributed Training of Large Language Models via Memory-Parallelism Co-Optimization

Zhanda Zhu^{1, 2, 3}, Christina Giannoula^{1, 2, 3}, Muralidhar Andoorveedu¹, Qidong Su^{1, 2, 3}, Karttikeya Mangalam⁴, Bojian Zheng^{1, 2, 3}, Gennady Pekhimenko^{1, 2, 3}









Executive Summary

- Accelerate distributed training by comprehensively co-optimizing <u>parallelism</u> and <u>memory optimizations</u>.
- Key Challenges
 - Exploded and complex configuration search space.
 - Inaccurate performance prediction: due to missing overlap and ignoring intermicrobatch imbalance.
- Mist addresses the challenges with
 - 1 Symbolic-based Performance Analysis
 - 2 Overlap-Centric Scheduling & Imbalance-aware hierarchical tuning
- Key Result: Up to $1.73 \times$ better vs. Megatron-LM and $2.04 \times$ vs. the automatic method Aceso, respectively.

Large Language Models

Widely adopted in both open communities and commercial systems.







• State-of-the-art performance in many applications:



Text Generation



Machine Translation

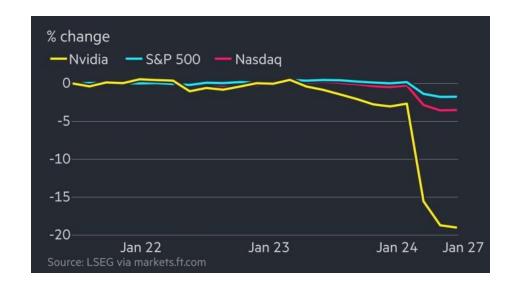


Speech Recognition

Cost of Training Large Language Models

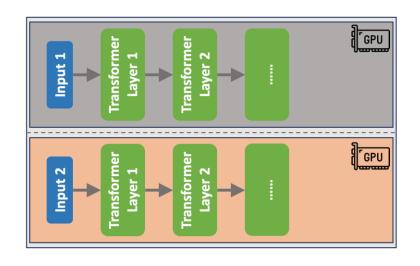
Extremely resource-consuming and expensive!

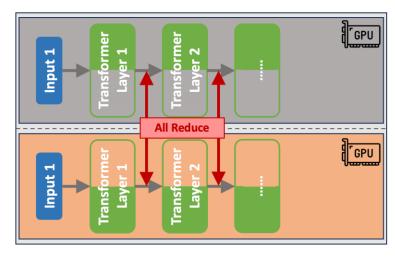
	Training Costs (\$)	Gas Emissions (tons CO ₂)
Llama 3.1 8B	2.92M	420
Llama 3.1 70B	14.00M	2,040
Llama 3.1 405B	61.68M	8,930
Total	78.60M	11,390

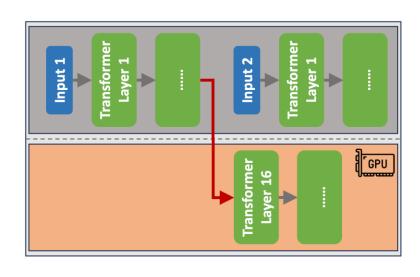


Strong incentive to reduce the training time of these models.

Optimizations – Distributed Parallelism







Data Parallelism

Partition input data

- Computation Utilization: 😉
- Communication Overhead: (2)
- Memory Usage: 😨

Tensor Parallelism

Partition linear layers

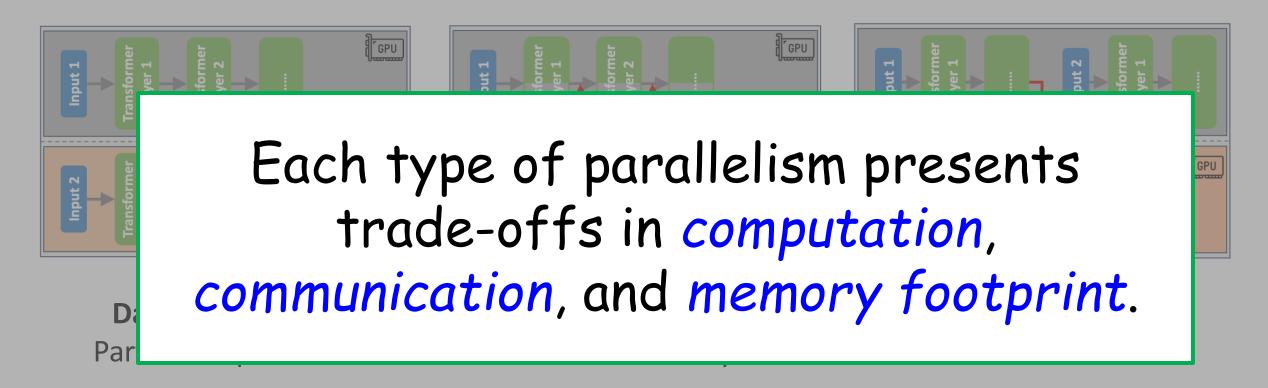
- Computation Overhead: (2)
 Communication Overhead: (2)
- Memory Usage: 😉

Pipeline Parallelism

Partitions the model into stages

- Computation Utilization: 💿
- Communication Overhead: (4)
- 🔻 Memory Usage: 😉

Optimizations – Distributed Parallelism

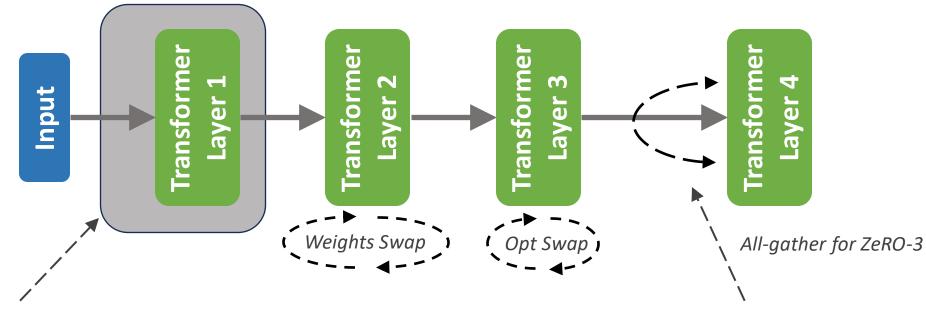


- Computation Utilization: (2)
- Communication Overhead: (2)
- Memory Usage:

- Computation Overhead: (2)
 Communication Overhead: (2)
- Memory Usage: (2)

- 🔻 Computation Utilization: 😧
- 🕨 Communication Overhead: 😉
- Memory Usage: 😀

Optimizations - Memory Optimizations



Activation Checkpointing (CKPT) Recompute activations in BWD

- Computation Utilization: (2)
 Communication Overhead: (2)
- Memory Usage: 💿

Offloading/Swapping
Offload tensors to CPU

- Computation Overhead: 😀 Communication Overhead: 😧
- Memory Usage: (2)

ZeRORedundancy elimination

- Computation Utilization:
- Communication Overhead: 😀
- 🕨 Memory Usage: 😀

Optimizations - Memory Optimizations



Each type of memory optimizations presents trade-offs in computation, communication, and memory footprint.

Activation Recomp

- Computation Utilization: (a)
- Communication Overhead: (4)
- Memory Usage:

- Computation Overhead: ©
 Communication Overhead: ©
- Memory Usage: (2)

- Computation Utilization: 😧
- Communication Overhead: 😉
- Memory Usage: (2)

Optimizations - Memory Optimizations



Parallelism and memory footprint reduction techniques need to be jointly optimized.

ActivationRecomp

- Computation Utilization: (4)
- Communication Overhead: (2)
- Memory Usage:

- Computation Overhead: (2)
 Communication Overhead: (2)
- Memory Usage: (2)

- Computation Utilization: ②
- Communication Overhead: (2)
- Memory Usage: (2)



Aggressive Memory Optimizations

Apply a higher level of ZeRO

Overhead ↑

Aggressive Memory Optimizations

- Apply a higher level of ZeRO
- Apply more CKPT

Overhead 个 个

Aggressive Memory Optimizations

- Apply a higher level of ZeRO
- Apply more CKPT
- Apply more offloading

Overhead 个 个 个

Aggressive Memory Optimizations

- Apply a higher level of ZeRO
- Apply more CKPT
- Apply more offloading

Overhead 个 个 个

Utilize Gained Memory

Reduce TP Size

Perf Gain 1

Aggressive Memory Optimizations

- Apply a higher level of ZeRO
- Apply more CKPT
- Apply more offloading

Overhead 个 个 个

Utilize Gained Memory

- Reduce TP Size
- Reduce PP Size

Perf Gain ↑ ↑

Aggressive Memory Optimizations

- Apply a higher level of ZeRO
- Apply more CKPT
- Apply more offloading

Overhead 个 个 个

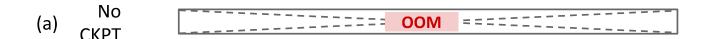


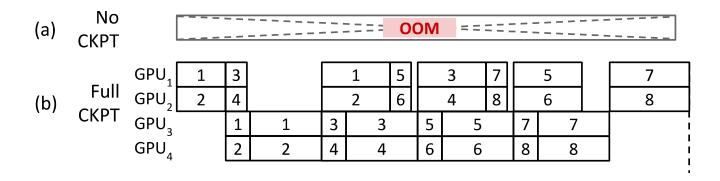
Utilize Gained Memory

- Reduce TP Size
- Reduce PP Size
- Increase batch size

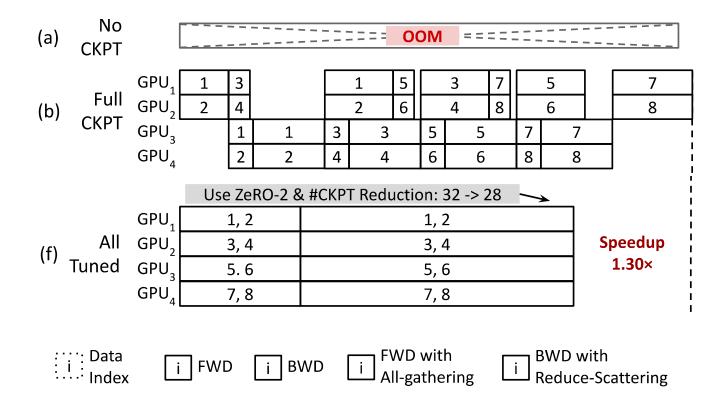




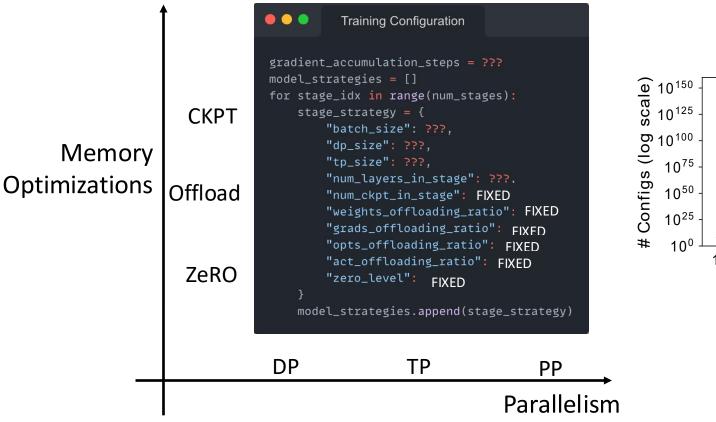


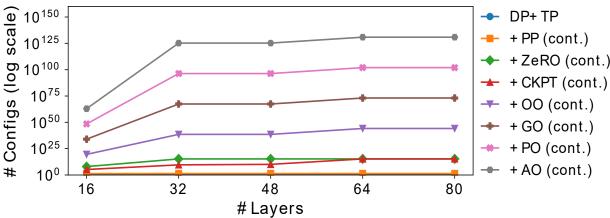






• Shortcoming #1: Unable to navigate the exploded search space.



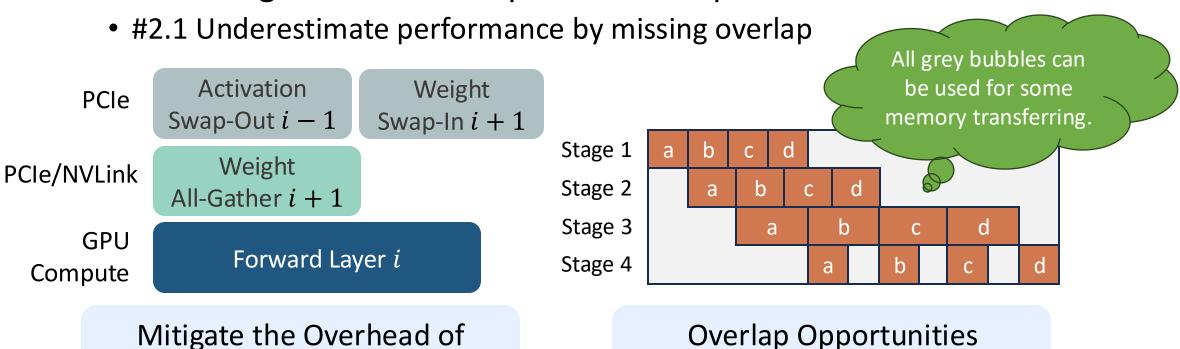


• Shortcoming #1: Unable to navigate the exploded search space.



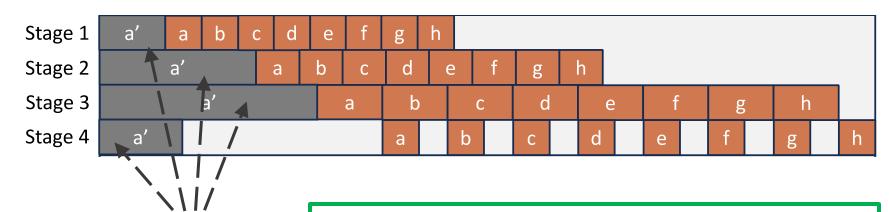
- Shortcoming #1: Unable to navigate the exploded search space.
- Shortcoming #2: Inaccurate performance prediction.

Offloading and ZeRO



in Pipeline Parallelism

- Shortcoming #1: Unable to navigate the exploded search space.
- Shortcoming #2: Inaccurate performance prediction.
 - #2.1 Under-estimate performance by missing overlap
 - #2.2 Mispredict performance by ignoring inter-microbatch imbalance



The first and last microbatches cost more time. extra (+-) ~7% error ratio for the inter-stage performance prediction, leading to up to extra 15% performance degradation.

Outline

Background

of distributed training and its optimizations

Shortcomings

of existing distributed training systems

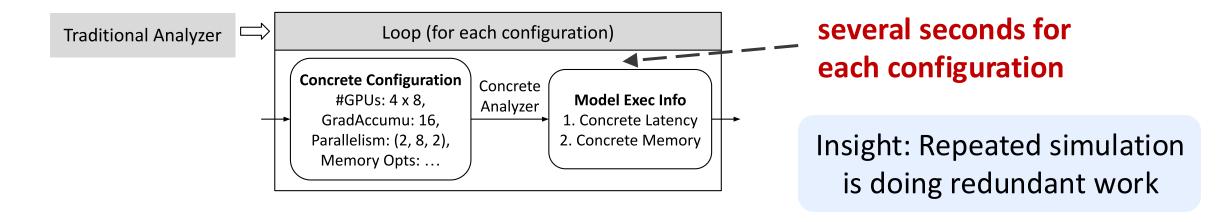
Key Ideas

of Mist to address these shortcomings

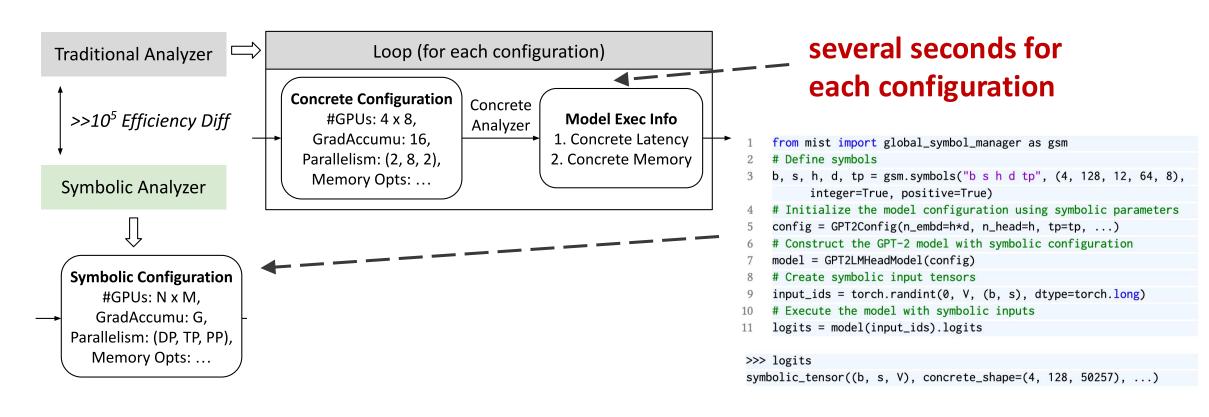
Evaluation

on state-of-the-art workloads

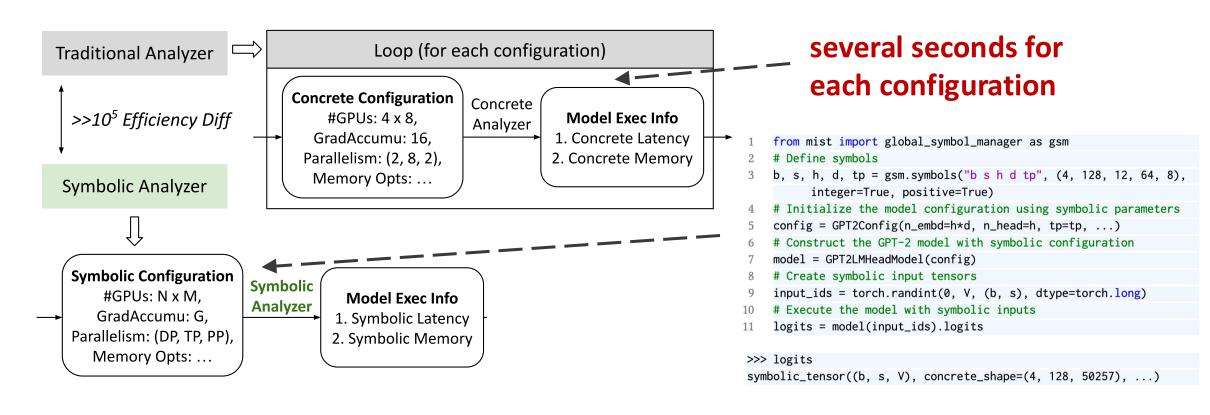
- Symbolic-Based Efficient Performance Prediction.
- 2) Fine-Grained Overlap-Centric Scheduling.
- (3) Imbalance-Aware Hierarchical Tuning via Pareto Frontier Sampling.
- 1 Symbolic-Based Efficient Performance Prediction



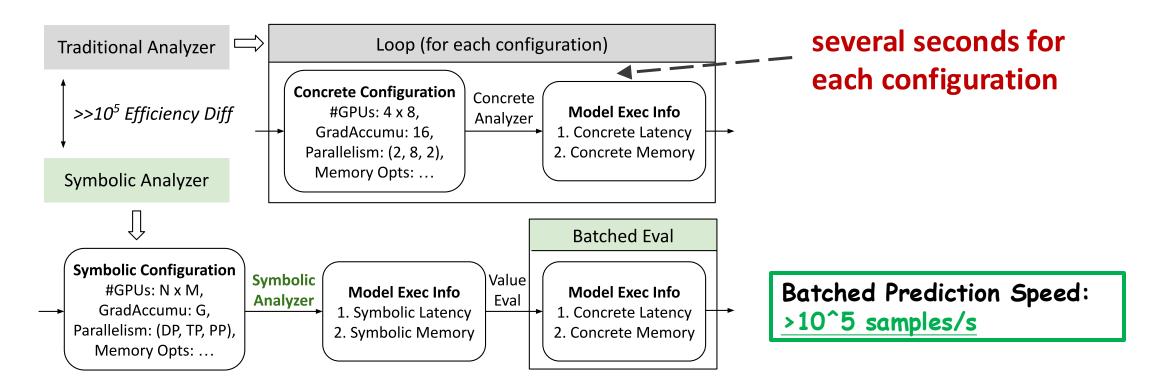
- 1 Symbolic-Based Efficient Performance Prediction.
- 2) Fine-Grained Overlap-Centric Scheduling.
- 3 Imbalance-Aware Hierarchical Tuning via Pareto Frontier Sampling.
- 1 Symbolic-Based Efficient Performance Prediction



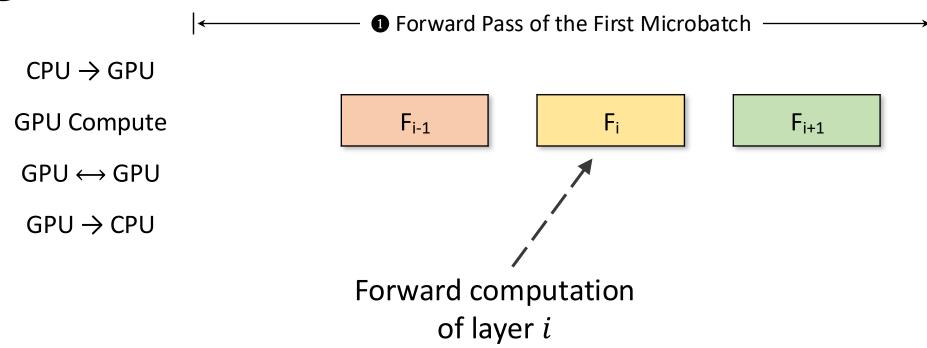
- 1 Symbolic-Based Efficient Performance Prediction.
- 2) Fine-Grained Overlap-Centric Scheduling.
- 3 Imbalance-Aware Hierarchical Tuning via Pareto Frontier Sampling.
- 1 Symbolic-Based Efficient Performance Prediction



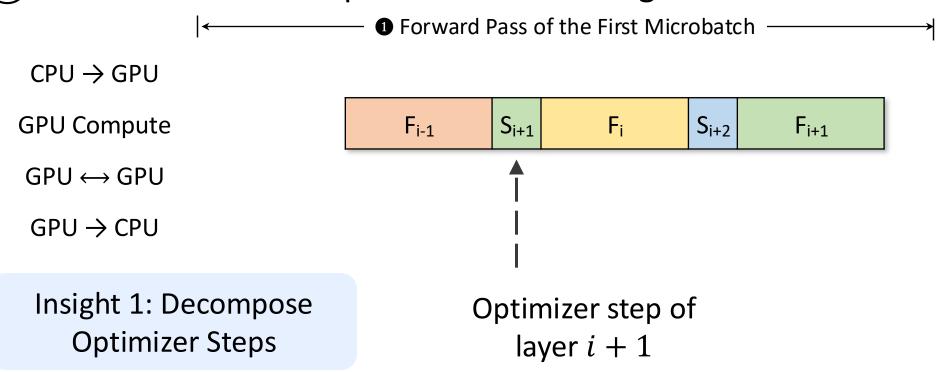
- 1 Symbolic-Based Efficient Performance Prediction.
- 2) Fine-Grained Overlap-Centric Scheduling.
- 3 Imbalance-Aware Hierarchical Tuning via Pareto Frontier Sampling.
- 1 Symbolic-Based Efficient Performance Prediction



- Symbolic-Based Efficient Performance Prediction.
- 2 Fine-Grained Overlap-Centric Scheduling.
- 3 Imbalance-Aware Hierarchical Tuning via Pareto Frontier Sampling.
- (2) Fine-Grained Overlap-Centric Scheduling.

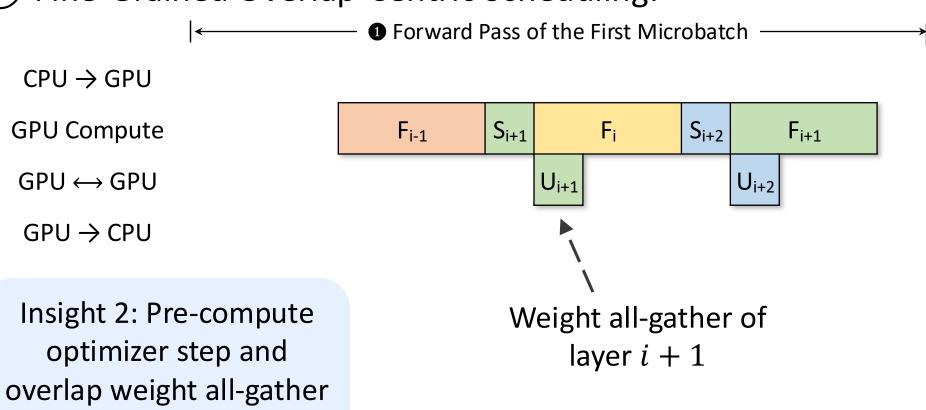


- Symbolic-Based Efficient Performance Prediction.
- 2 Fine-Grained Overlap-Centric Scheduling.
- (3) Imbalance-Aware Hierarchical Tuning via Pareto Frontier Sampling.
- (2) Fine-Grained Overlap-Centric Scheduling.

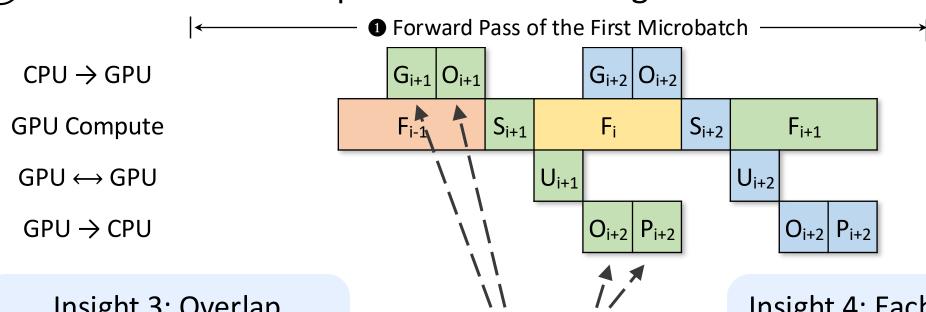


with computation

- 1) Symbolic-Based Efficient Performance Prediction.
- 2) Fine-Grained Overlap-Centric Scheduling.
- (3) Imbalance-Aware Hierarchical Tuning via Pareto Frontier Sampling.
- (2) Fine-Grained Overlap-Centric Scheduling.



- 1) Symbolic-Based Efficient Performance Prediction.
- 2 Fine-Grained Overlap-Centric Scheduling.
- 3 Imbalance-Aware Hierarchical Tuning via Pareto Frontier Sampling.
- (2) Fine-Grained Overlap-Centric Scheduling.

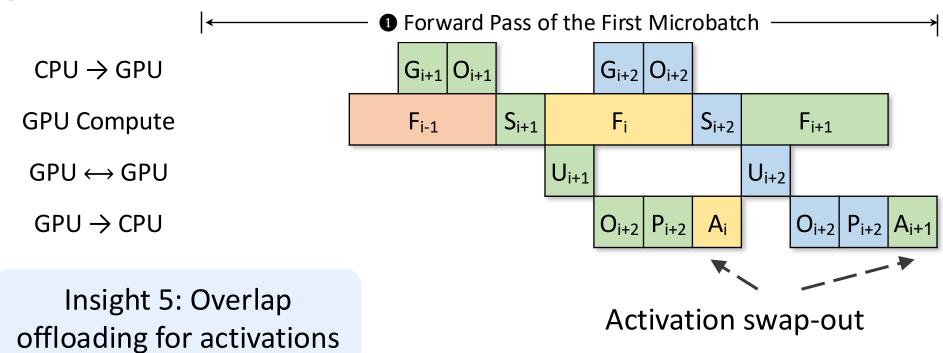


Insight 3: Overlap offloading for optimizer step related tensors

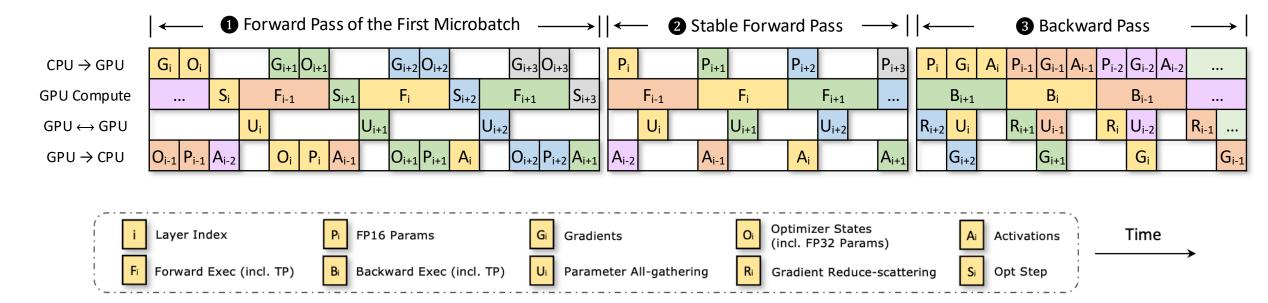
Gradients swap-in
Opt states swap-in
Opt states swap-out
Weights swap-out

Insight 4: Each offloading type can be set with a ratio from 0 to 1.

- Symbolic-Based Efficient Performance Prediction.
- 2 Fine-Grained Overlap-Centric Scheduling.
- 3 Imbalance-Aware Hierarchical Tuning via Pareto Frontier Sampling.
- (2) Fine-Grained Overlap-Centric Scheduling.

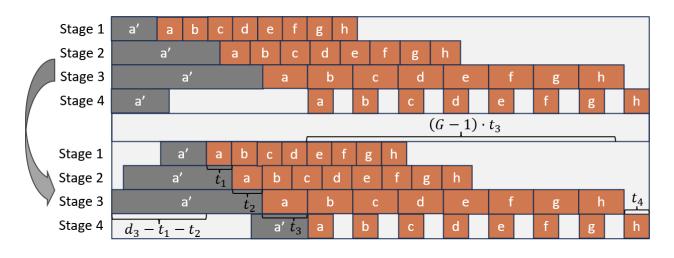


- 1) Symbolic-Based Efficient Performance Prediction.
- 2 Fine-Grained Overlap-Centric Scheduling.
- 3 Imbalance-Aware Hierarchical Tuning via Pareto Frontier Sampling.
- 2 Fine-Grained Overlap-Centric Scheduling.



- 1) Symbolic-Based Efficient Performance Prediction.
- 2) Fine-Grained Overlap-Centric Scheduling.
- (3) Imbalance-Aware Hierarchical Tuning via Pareto Frontier Sampling.
- (3) Imbalance-Aware Hierarchical Tuning via Pareto Frontier Sampling.
 - Formulize inter-stage tuning as a MILP problem.

However, t_i and d_i are correlated.



Inter-Stage $\min_{\substack{1 \leq i \leq S \\ l_i, f_i, (n_i, m_i)}} \left\{ (G-1) \cdot \max_{1 \leq i \leq S} \left\{ t_i \right\} + \sum_{i=1}^S t_i + \max_{1 \leq i \leq S} \left(d_i - \sum_{1 \leq j < i} t_i \right) \right\}$ Normal pipeline time calculation microbatch differences

Intra-Stage Sample (t_i, d_i) from the pareto frontier of intra stage $\min_{p,z,o} \alpha \cdot G \cdot t_{p,z,o} + (1-\alpha) \cdot d_{p,z,o}$ $i.e. \max \left(Mem_{peak}^{(fwd)}, Mem_{peak}^{(bwd)} \right) \leq Mem_{Budget}$

 t_i : Time of stable micro batch in stage i

 d_i : Time difference between the first microbatch and t_i in stage i

Mist Overview

An automatic distributed training performance analysis and tuning system with overlapped execution engine.

OverlapCentric
Schedule
Template
(§4.1)

Outline

Background

of distributed training and its optimizations

Shortcomings

of existing distributed training systems

Key Ideas

of Mist to address these shortcomings

Evaluation

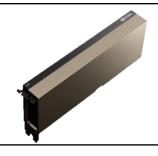
on state-of-the-art workloads

Methodology





NVIDIA L4 GPU



NVIDIA A100 GPU

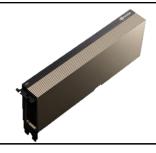
Platform	GPU	GPU#	Mem.	PCIe Spec	NVLink	Interconnect
GCP AWS		[2, 4, 8, 16, 32] [2, 4, 8, 16, 32]		_		100Gbps 400Gbps

Methodology





NVIDIA L4 GPU

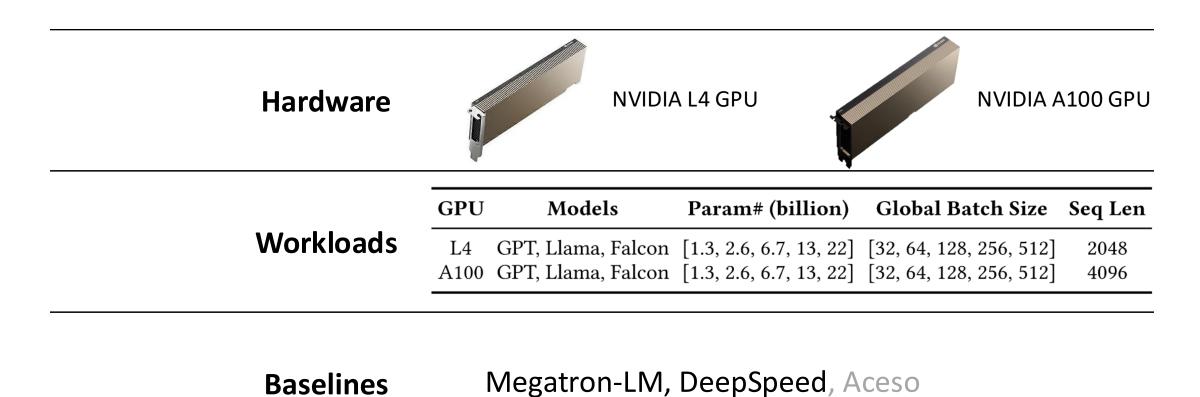


NVIDIA A100 GPU

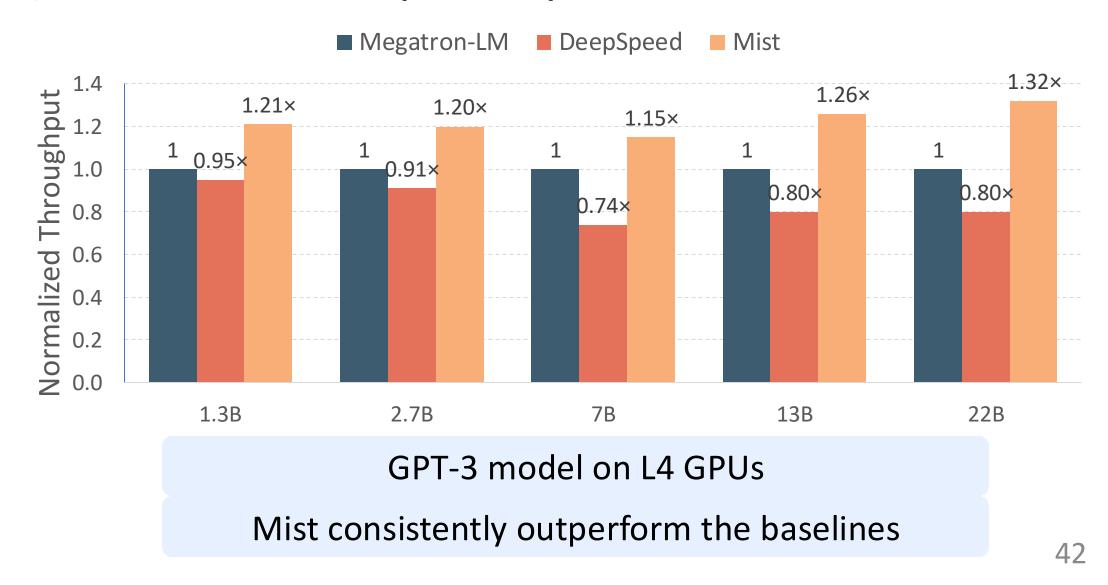
Workloads

GPU	Models	Param# (billion)	Global Batch Size	Seq Len
			[32, 64, 128, 256, 512] [32, 64, 128, 256, 512]	2048 4096

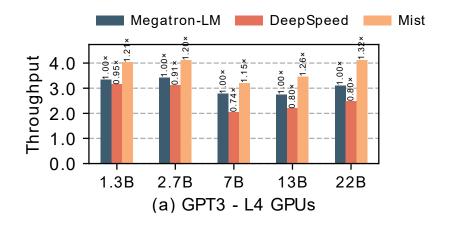
Methodology



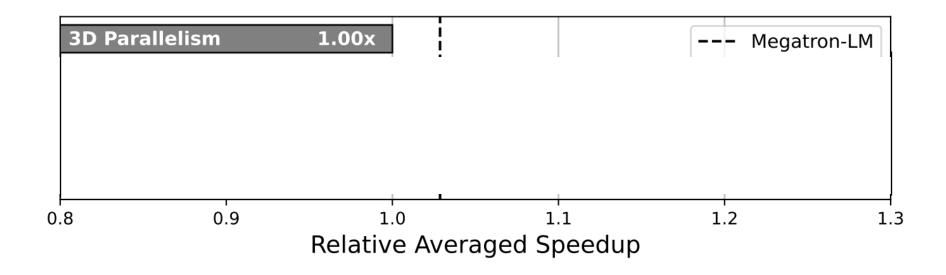
1) Performance Speedup w/ FlashAttn

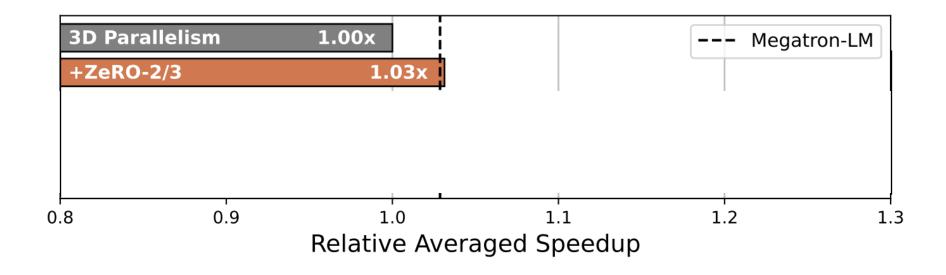


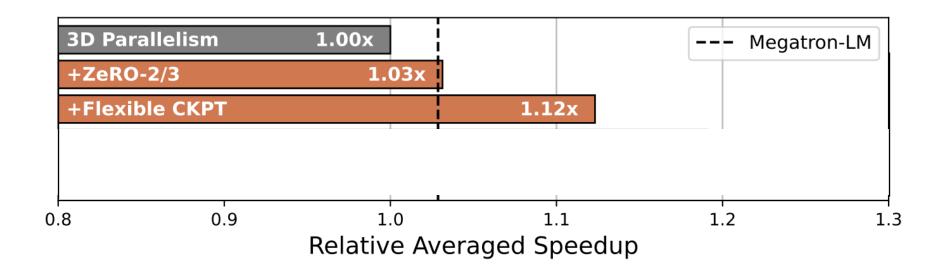
1) Performance Speedup w/ FlashAttn

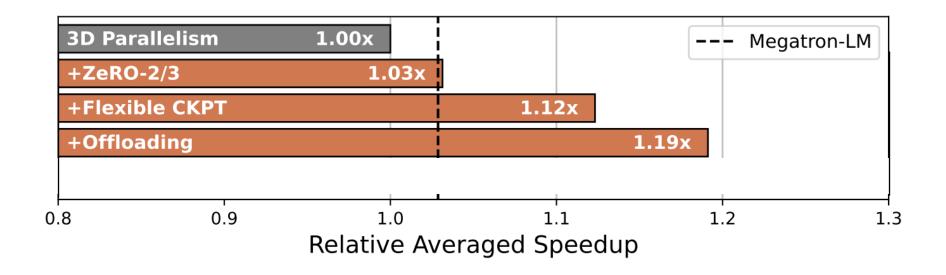


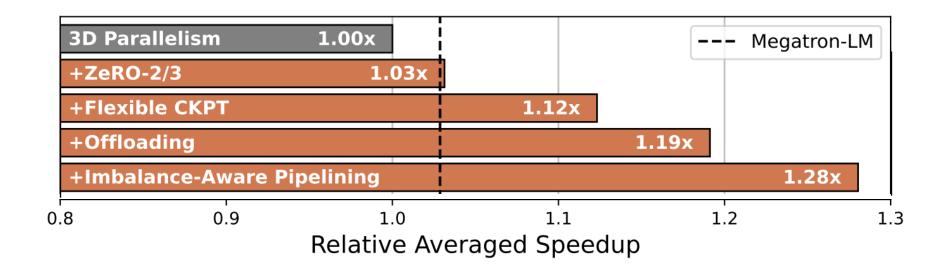
- L4 GPUs: $1.32 \times$ on average (up to $1.59 \times$) over Megatron-LM
- A100 GPUs: $1.34 \times$ on average (up to $1.72 \times$) over Megatron-LM



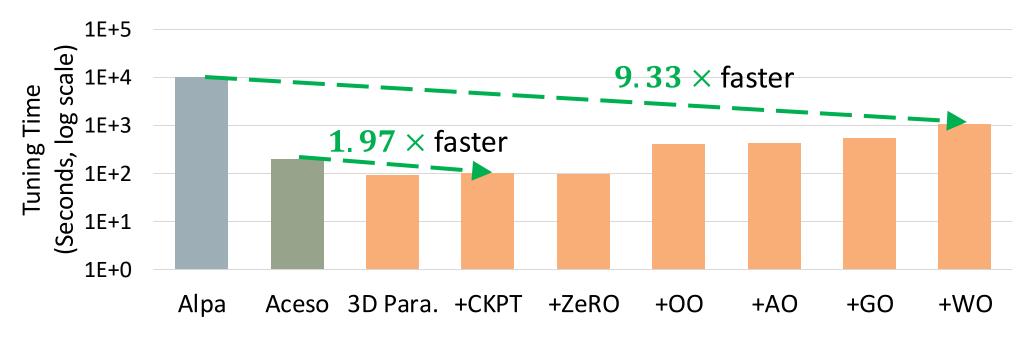








3) Tuning Time



- Vs. Aceso: with the same search space,
 Mist is 1. 97 × faster.
- Vs. Alpa: even with larger search space,
 Mist is 9. 33 × faster.

Executive Summary







- Mist accelerates distributed training by co-optimizing <u>parallelism</u> and <u>memory optimizations</u>.
- Key Challenges
 - Exploded and complex configuration search space.
 - Inaccurate performance prediction: due to missing overlap and ignoring intermicrobatch imbalance.
- Mist addresses the challenges with
 - 1 Symbolic-based Performance Analysis
 - 2 Overlap-Centric Scheduling & Imbalance-aware hierarchical tuning
- Key Result: Up to $1.73 \times$ better vs. Megatron-LM and $2.04 \times$ vs. the automatic method Aceso, respectively.