

Accelerating Distributed MoE Training and Inference with Lina

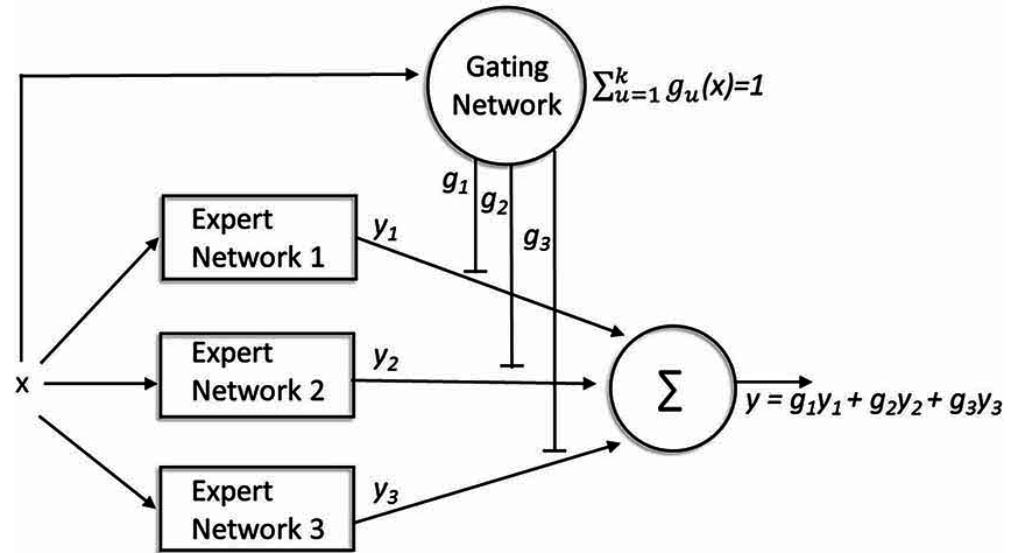
Jiamin Li¹, Yimin Jiang², Yibo Zhu, Cong Wang¹, Hong Xu²

¹City University of Hong Kong, ²ByteDance Inc., ³The Chinese University of Hong Kong

ATC 2023

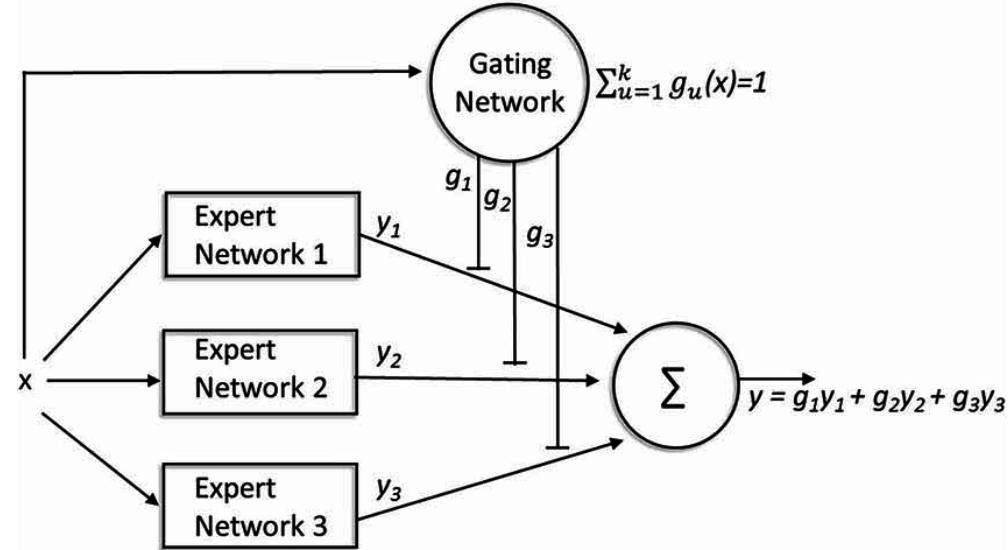


Sparsely-Activated Mixture-of-Experts (MoE)



MoE architecture
An ensemble of experts.

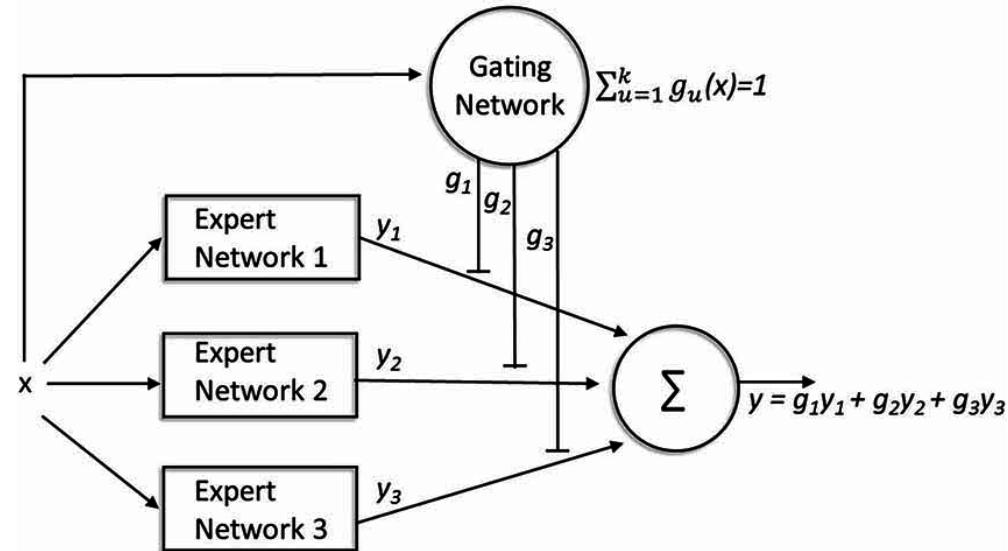
Sparsely-Activated Mixture-of-Experts (MoE)



MoE architecture
An ensemble of experts.

- Sparsely-activated MoE: each input selects just a few (1 or 2) experts for processing
- Benefit: sub-linear scaling of FLOPS with model size

Sparsely-Activated Mixture-of-Experts (MoE)



MoE architecture
An ensemble of experts.

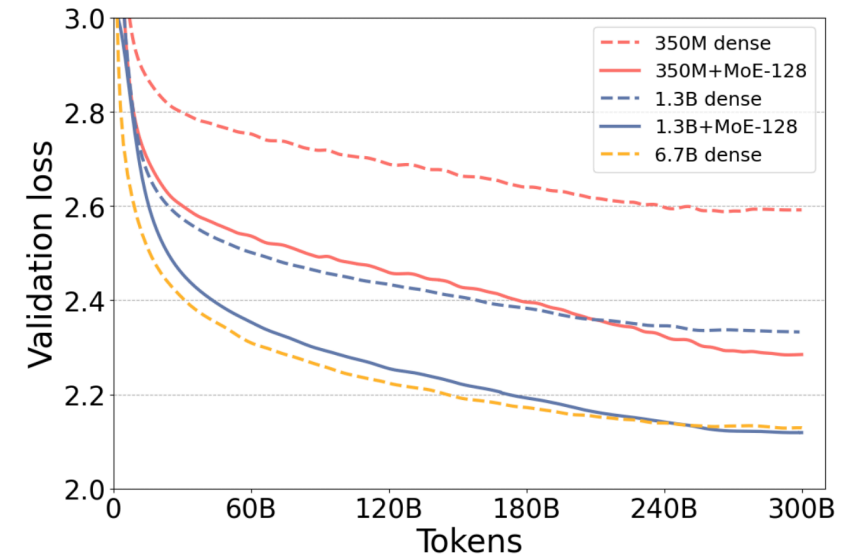
- Sparsely-activated MoE: each input selects just a few (1 or 2) experts for processing
- Benefit: sub-linear scaling of FLOPS with model size

Massive model parameters with constant computation cost.

Potential of MoE in Transformer Models

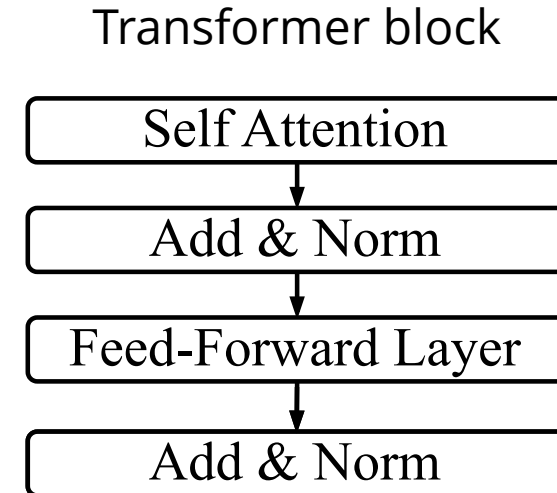
- GLaM by Google
 - GLaM outperforms GPT-3 on 29 tasks
- DeepSpeed MoE models
 - Model quality: 6.7B-parameter dense = 1.3B-parameter MoE - 128
 - Training compute reduction of 5x

		GPT-3	GLaM	relative
cost	FLOPs / token (G)	350	180	-48.6%
	Train energy (MWh)	1287	456	-64.6%
accuracy on average	Zero-shot	56.9	62.7	+10.2%
	One-shot	61.6	65.5	+6.3%
	Few-shot	65.2	68.1	+4.4%



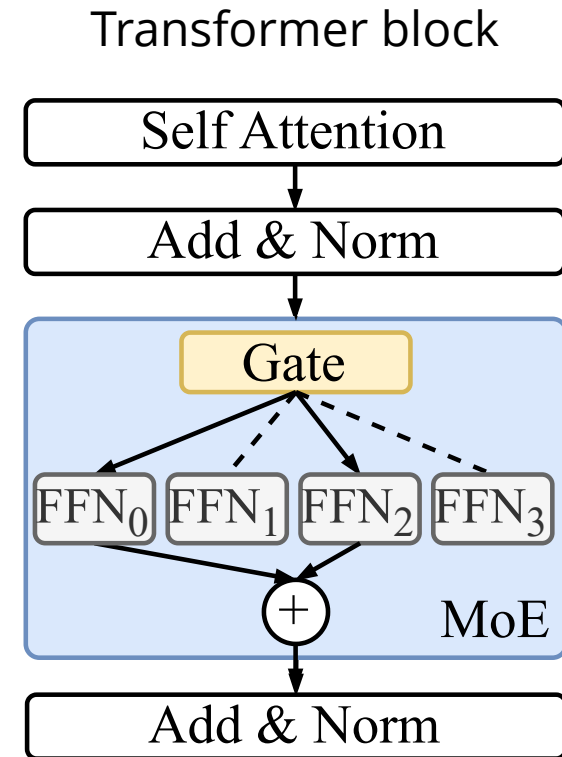
MoE in Transformer-based Language Models

- Feed forward layers (FFN) are replaced with MoE layers.
- MoE layer = gate + experts
 - Expert: feed forward neural network
=> same architecture, different parameters
 - Gate: a trainable matrix to select expert for each data sample
 - Top-2 in training, Top-1 in inference



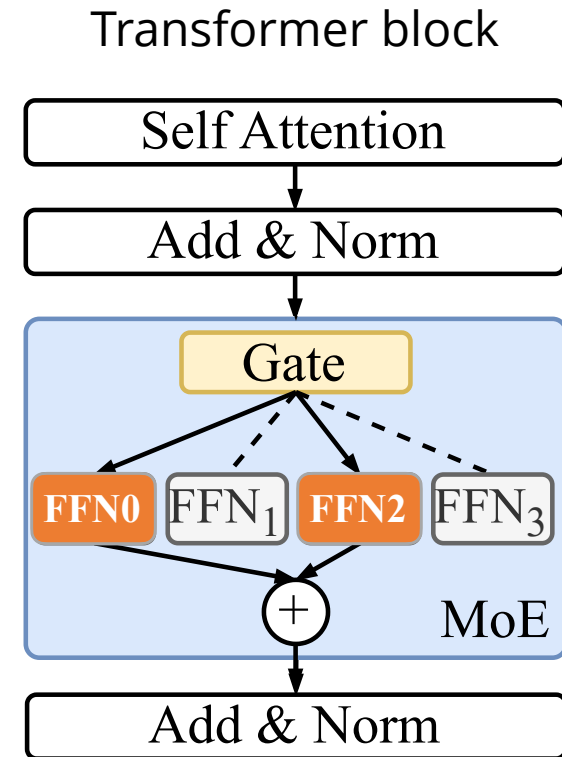
MoE in Transformer-based Language Models

- Feed forward layers (FFN) are replaced with MoE layers.
- MoE layer = gate + experts
 - Expert: feed forward neural network
=> same architecture, different parameters
- Gate: a trainable matrix to select expert for each data sample
 - Top-2 in training, Top-1 in inference



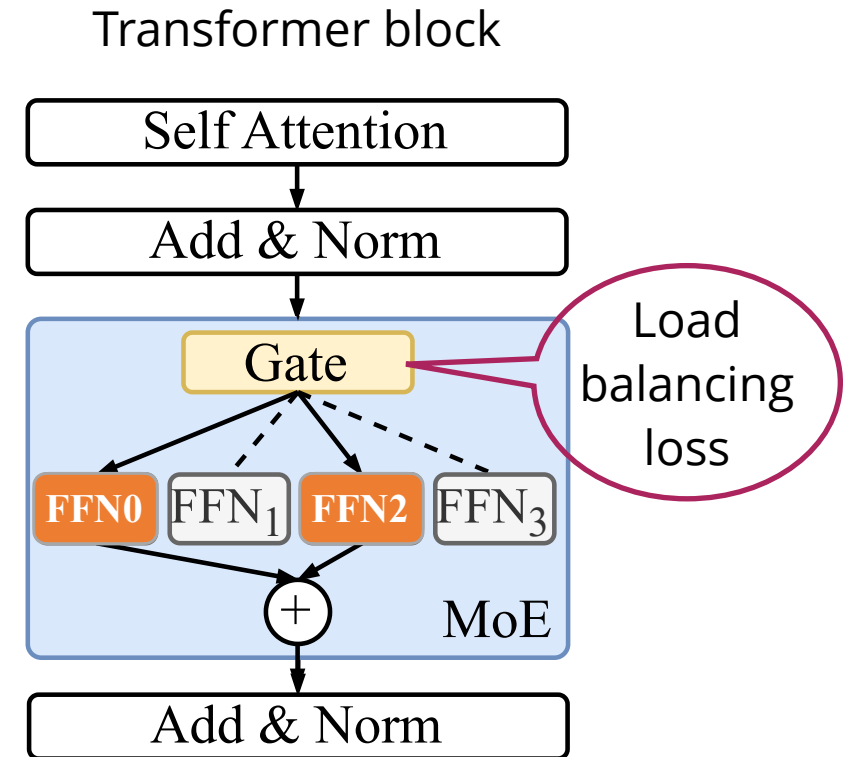
MoE in Transformer-based Language Models

- Feed forward layers (FFN) are replaced with MoE layers.
- MoE layer = gate + experts
 - Expert: feed forward neural network
=> same architecture, different parameters
- Gate: a trainable matrix to select expert for each data sample
 - Top-2 in training, Top-1 in inference



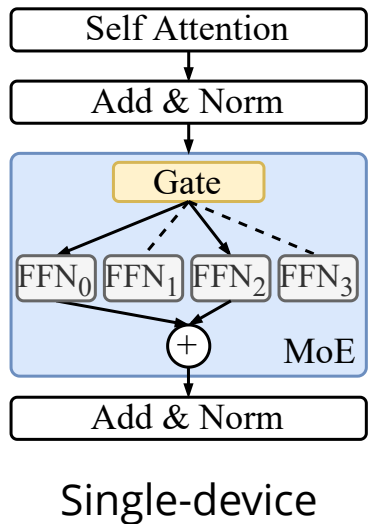
MoE in Transformer-based Language Models

- Feed forward layers (FFN) are replaced with MoE layers.
- MoE layer = gate + experts
 - Expert: feed forward neural network
=> same architecture, different parameters
 - Gate: a trainable matrix to select expert for each data sample
 - Top-2 in training, Top-1 in inference
 - Load balancing loss during training: even distribution among experts.



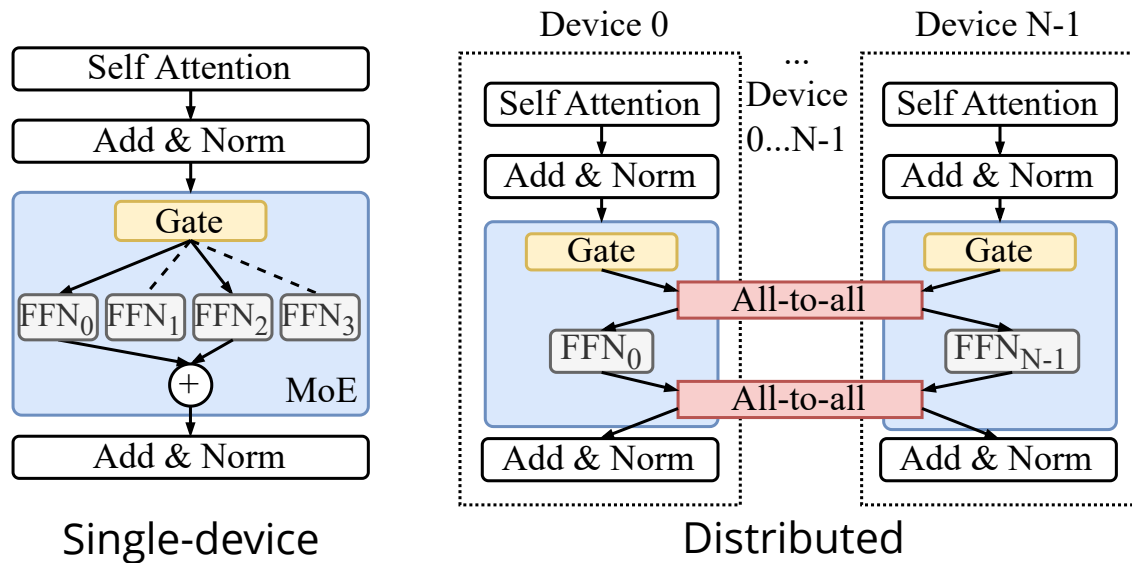
Distributed MoE

- Hybrid parallelism:
 - Expert parallelism: each device hosts one unique expert
 - Data parallelism: replicate non-expert parameter on each device



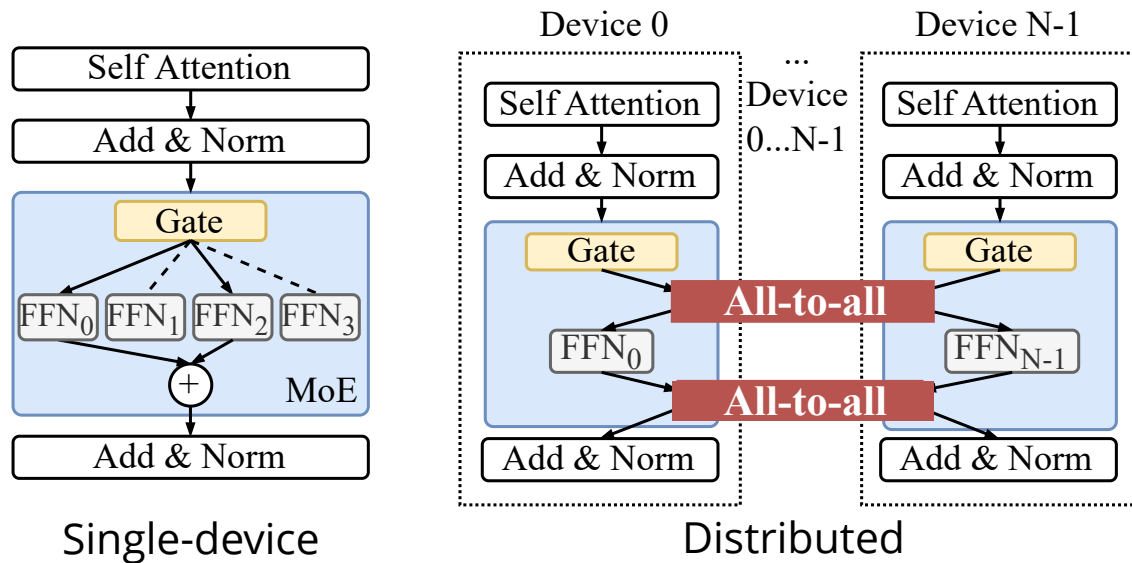
Distributed MoE

- Hybrid parallelism:
 - Expert parallelism: each device hosts one unique expert
 - Data parallelism: replicate non-expert parameter on each device



Distributed MoE

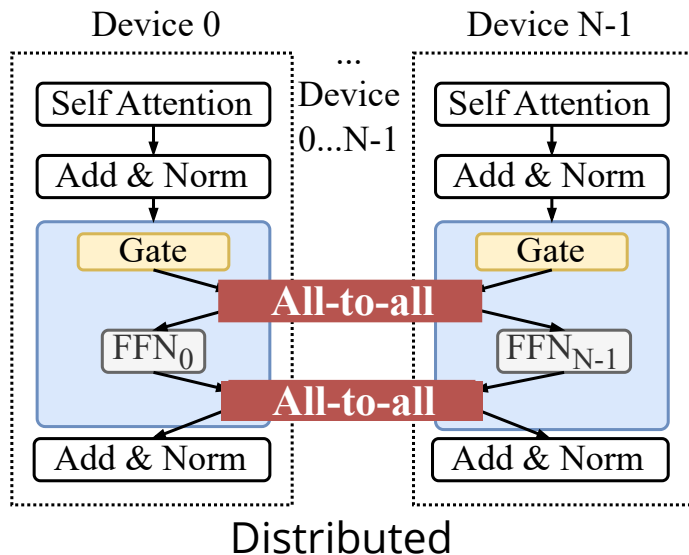
- Hybrid parallelism:
 - Expert parallelism: each device hosts one unique expert
 - Data parallelism: replicate non-expert parameter on each device



- All-to-all communication
 - 1st: send data samples to experts.
 - 2nd: restore data samples back

Distributed MoE

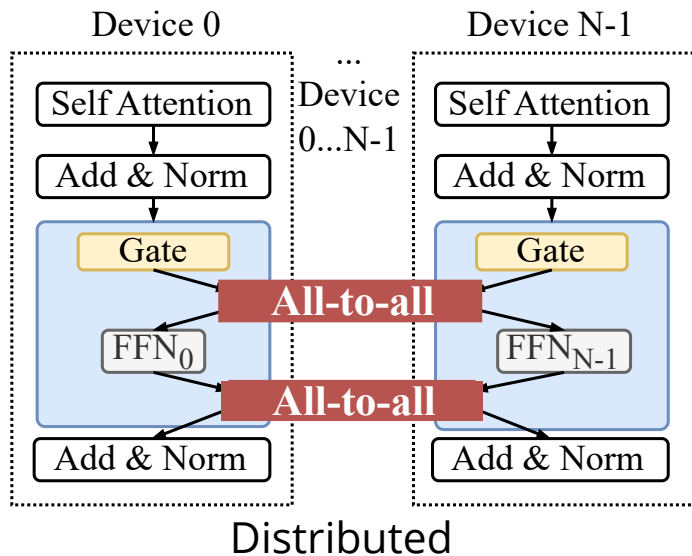
- Hybrid parallelism:
 - Expert parallelism: each device hosts one unique expert
 - Data parallelism: replicate non-expert parameter on each device



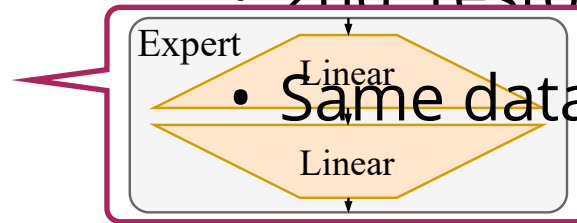
- All-to-all communication
 - 1st: send data samples to experts.
 - 2nd: restore data samples back

Distributed MoE

- Hybrid parallelism:
 - Expert parallelism: each device hosts one unique expert
 - Data parallelism: replicate non-expert parameter on each device



- All-to-all communication
 - 1st: send data samples to experts.
 - 2nd: restore data samples back
 - Same data transfer size



Distributed MoE is not efficient

Model #Layers & Params	Training (ms)		Inference (ms)	
	All-to-all	Ratio	All-to-all	Ratio
12L + 117M	259	36.7%	73	27.4%
24L + 233M	589	35.4%	103	26.2%
36L + 349M	979	38.2%	153	28.3%
12L + 419M	333	39.5%	102	32.5%
24L + 838M	715	37.6%	177	31.7%
36L + 1.2B	1145	36.8%	243	27.4%

All-to-all takes an average of 34.1% of the step time.

Distributed MoE is not efficient

All-to-all takes an average of 34.1% of the step time.

Time (ms)

Ratio

27.4%

26.2%

28.3%

32.5%

31.7%

27.4%

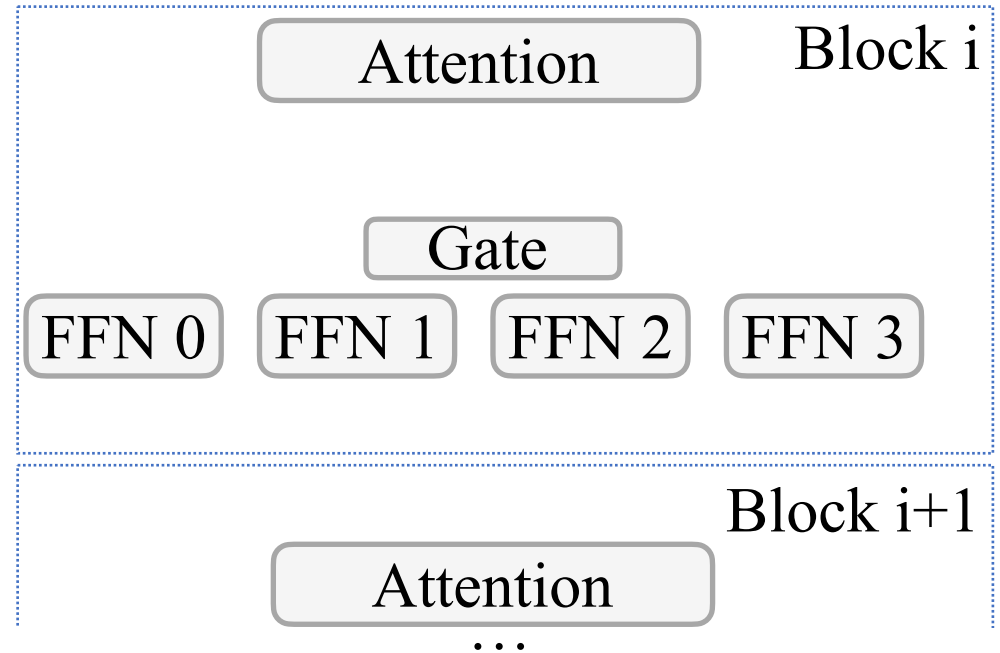
Distributed MoE is not efficient

All-to-all takes an average of 34.1% of the step time.

Component	Ratio
Attention	27.4%
Gate	26.2%
FFN 0	28.3%
FFN 1	32.5%
FFN 2	31.7%
FFN 3	27.4%

In MoE Transformer:

- Attention -> a complete sequence
- FFN expert -> one single token



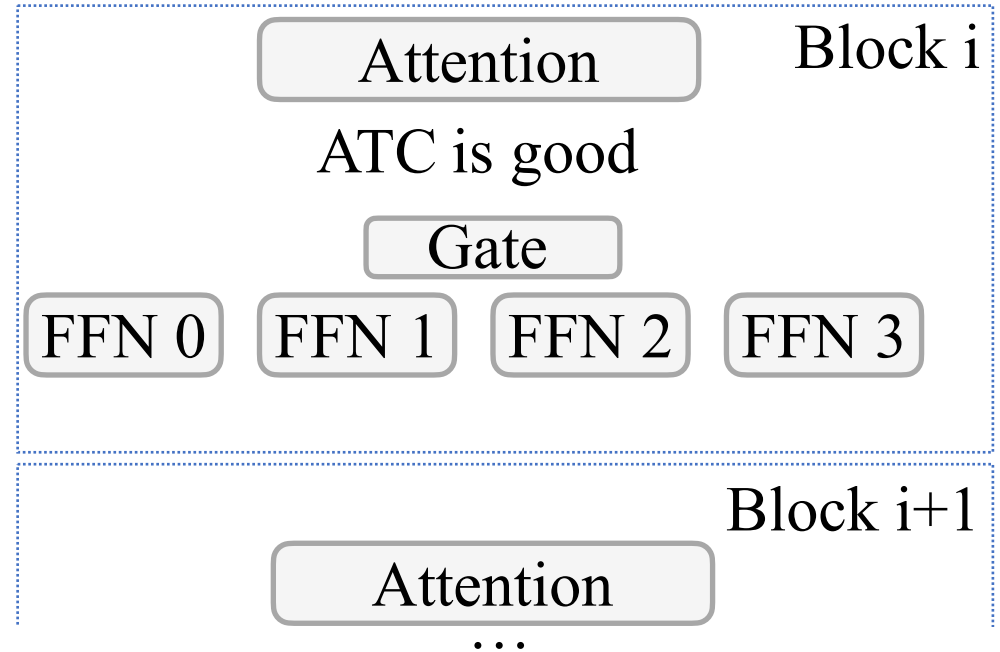
Distributed MoE is not efficient

All-to-all takes an average of 34.1% of the step time.

Component	Ratio
Attention	27.4%
FFN 0	26.2%
FFN 1	28.3%
FFN 2	32.5%
FFN 3	31.7%
Gate	27.4%

In MoE Transformer:

- Attention -> a complete sequence
- FFN expert -> one single token



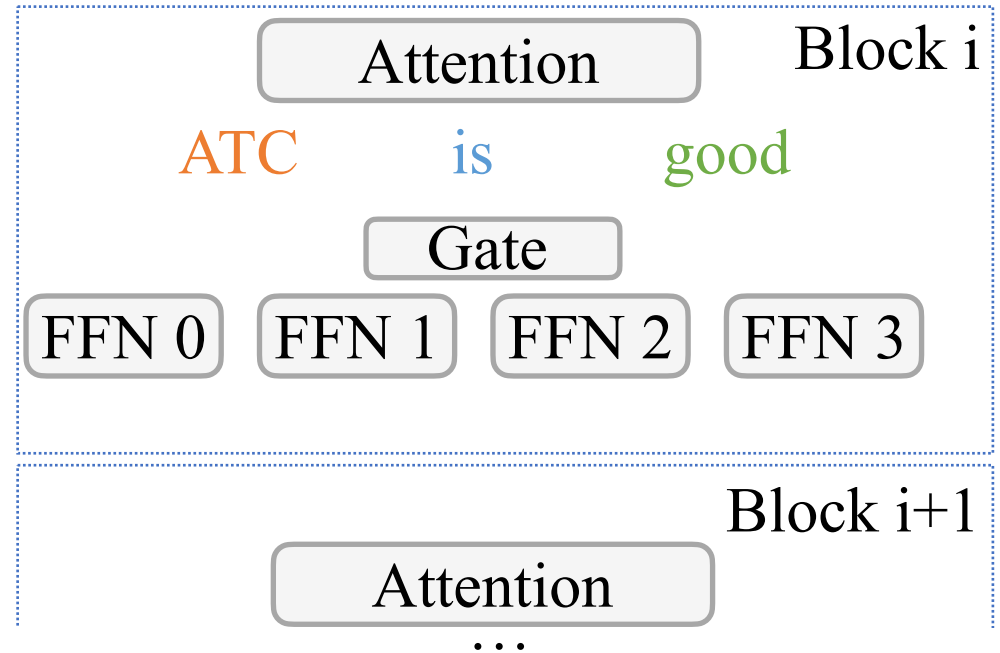
Distributed MoE is not efficient

All-to-all takes an average of 34.1% of the step time.

Time (ms)	Ratio
27.4%	
26.2%	
28.3%	
32.5%	
31.7%	
27.4%	

In MoE Transformer:

- Attention -> a complete sequence
- FFN expert -> one single token



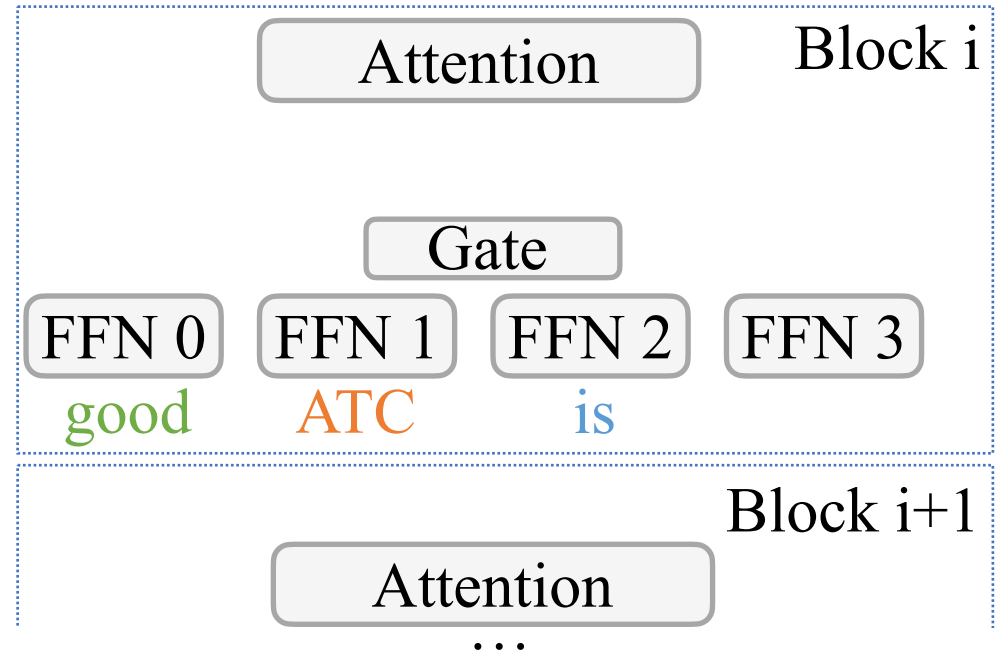
Distributed MoE is not efficient

All-to-all takes an average of 34.1% of the step time.

Time (ms)	Ratio
27.4%	
26.2%	
28.3%	
32.5%	
31.7%	
27.4%	

In MoE Transformer:

- Attention -> a complete sequence
- FFN expert -> one single token



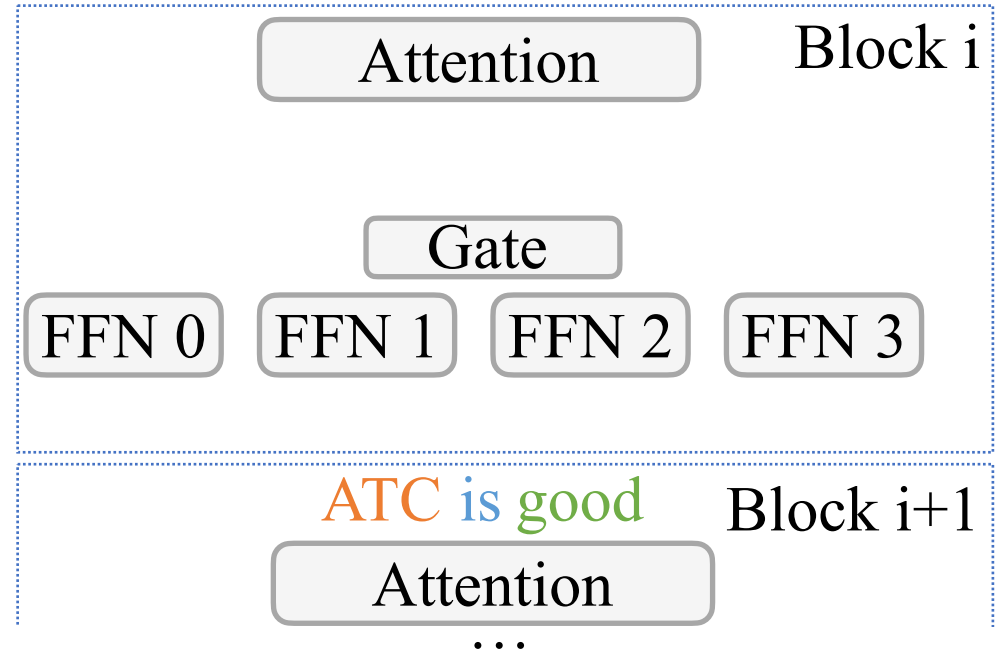
Distributed MoE is not efficient

All-to-all takes an average of 34.1% of the step time.

Component	Ratio
Attention	27.4%
Gate	26.2%
FFN 0	28.3%
FFN 1	32.5%
FFN 2	31.7%
FFN 3	27.4%

In MoE Transformer:

- Attention -> a complete sequence
- FFN expert -> one single token



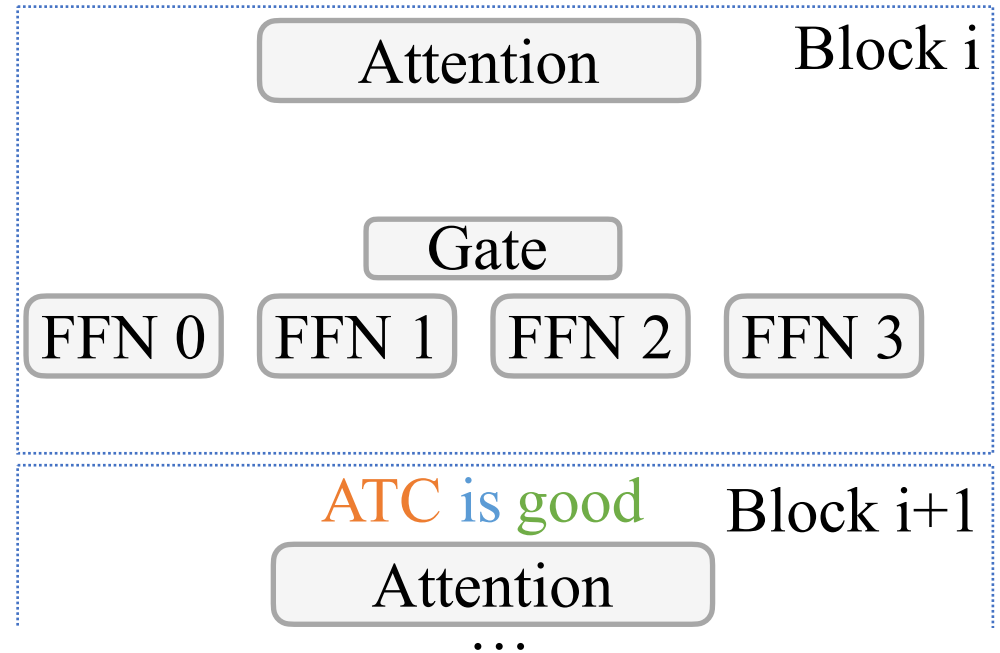
Distributed MoE is not efficient

All-to-all takes an average of 34.1% of the step time.

Component	Ratio
Attention	27.4%
Gate	26.2%
FFN 0	28.3%
FFN 1	32.5%
FFN 2	31.7%
FFN 3	27.4%

In MoE Transformer:

- Attention -> a complete sequence
- FFN expert -> one single token



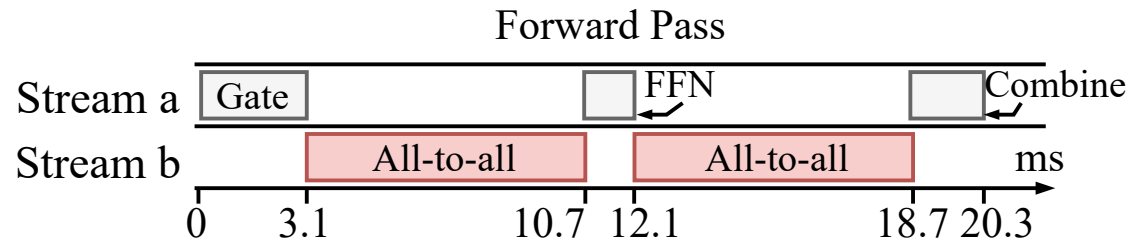
To restore to a sequence, we must wait for the processing of all tokens.

All-to-all is the bottleneck

- Synchronous and blocking operation & large amounts of data transfer.

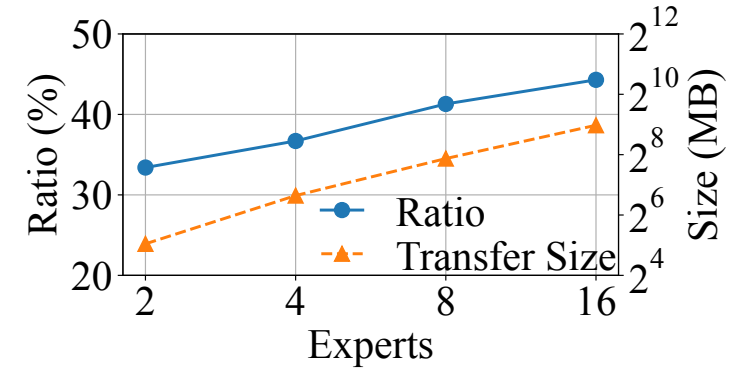
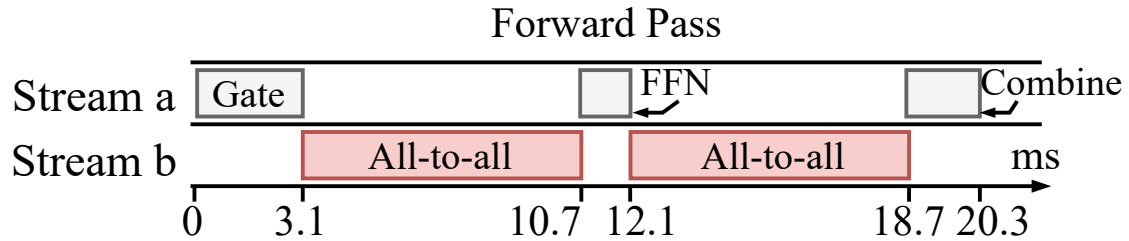
All-to-all is the bottleneck

- Synchronous and blocking operation & large amounts of data transfer.



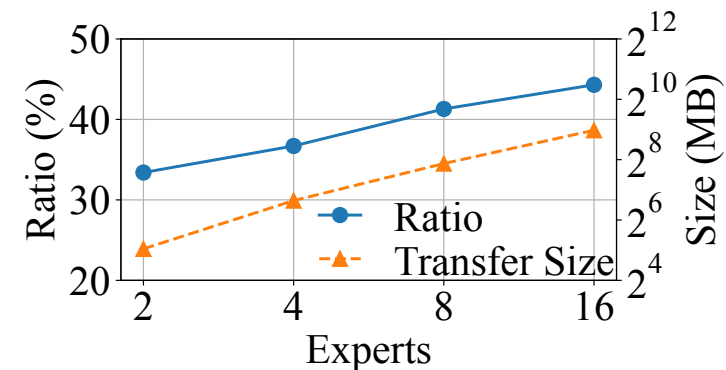
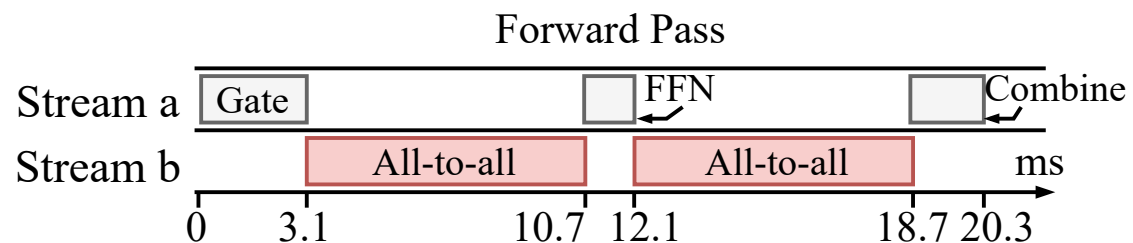
All-to-all is the bottleneck

- Synchronous and blocking operation & large amounts of data transfer.



All-to-all is the bottleneck

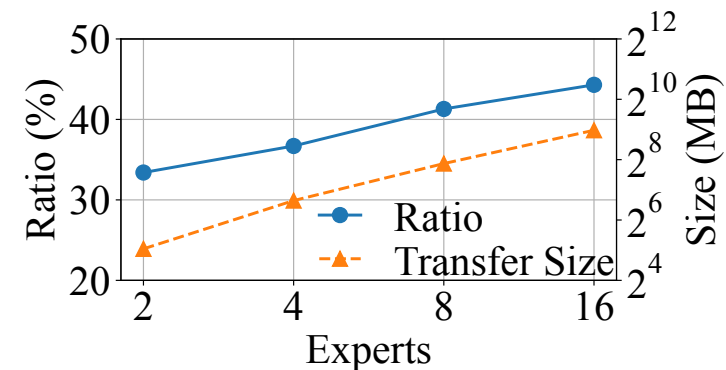
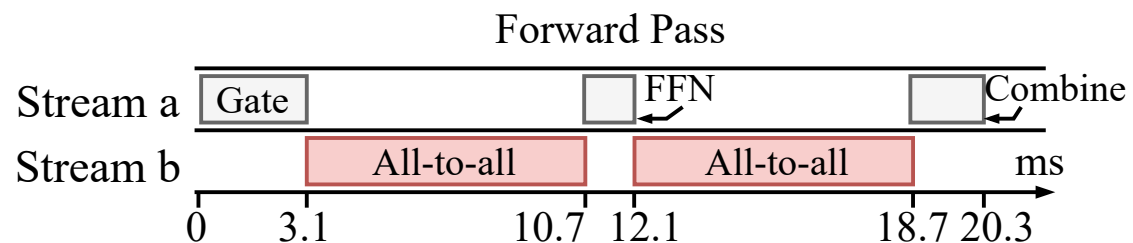
- Synchronous and blocking operation & large amounts of data transfer.



Is it the only cause of all-to-all being the bottleneck?

All-to-all is the bottleneck

- Synchronous and blocking operation & large amounts of data transfer.

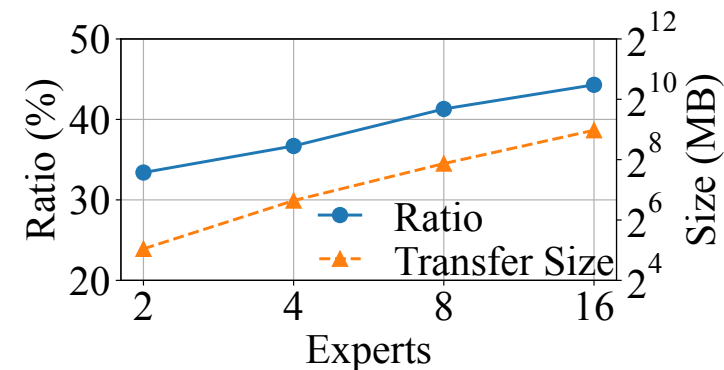
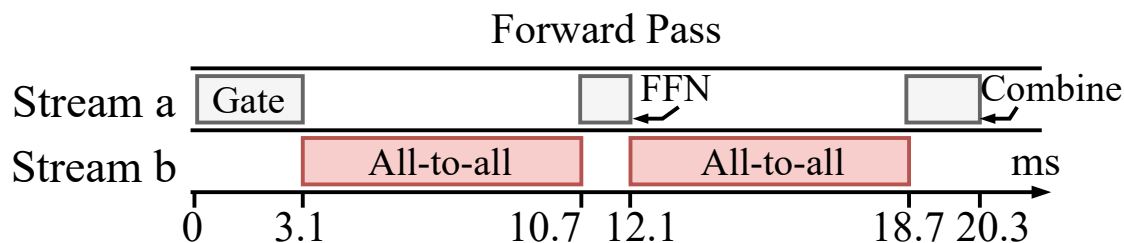


Is it the only cause of all-to-all being the bottleneck?

- MoE training and inference have their *unique* problems.

All-to-all is the bottleneck

- Synchronous and blocking operation & large amounts of data transfer.



Is it the only cause of all-to-all being the bottleneck?

- MoE training and inference have their *unique* problems.
 - Training has backward pass
 - Inference is purely workload-driven

MoE Training in Data & Expert Parallel

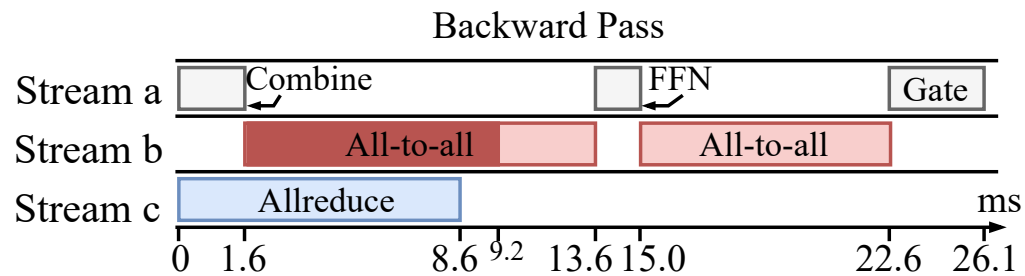
In backward pass,

- Allreduce: asynchronously aggregate non-expert gradients in data parallel.
- All-to-all: exchange token gradients to compute expert gradients.

MoE Training in Data & Expert Parallel

In backward pass,

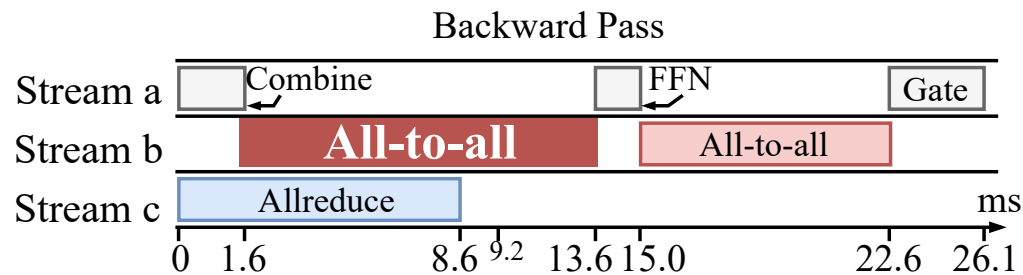
- Allreduce: asynchronously aggregate *non-expert* gradients in data parallel.
- All-to-all: exchange token gradients to compute *expert* gradients.



MoE Training in Data & Expert Parallel

In backward pass,

- Allreduce: asynchronously aggregate *non-expert* gradients in data parallel.
- All-to-all: exchange token gradients to compute *expert* gradients.

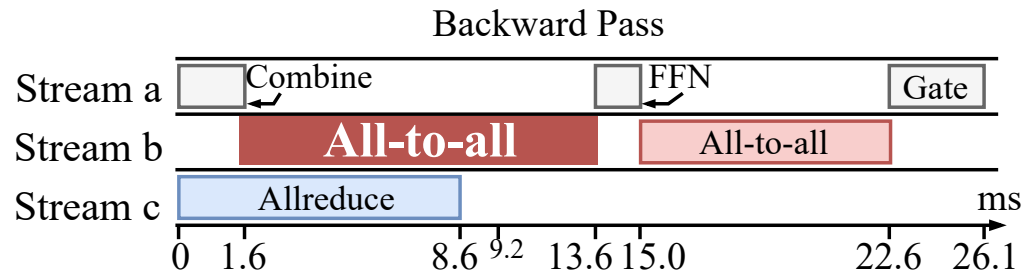


All-to-all is prolonged when it overlaps with allreduce and directly impacts step time.

MoE Training in Data & Expert Parallel

In backward pass,

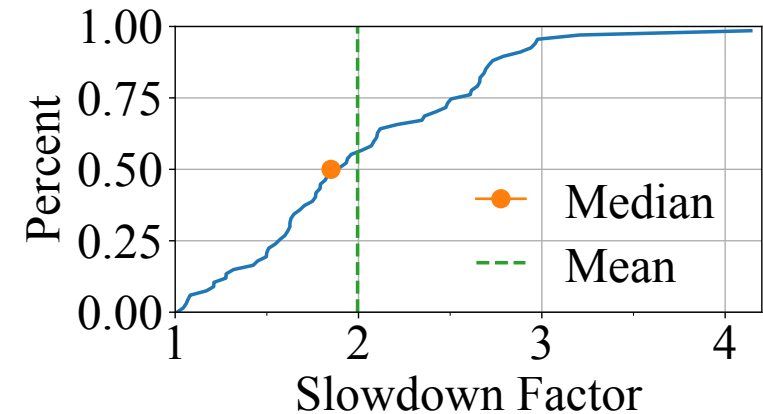
- Allreduce: asynchronously aggregate *non-expert* gradients in data parallel.
- All-to-all: exchange token gradients to compute *expert* gradients.



All-to-all is prolonged when it overlaps with allreduce and directly impacts step time.

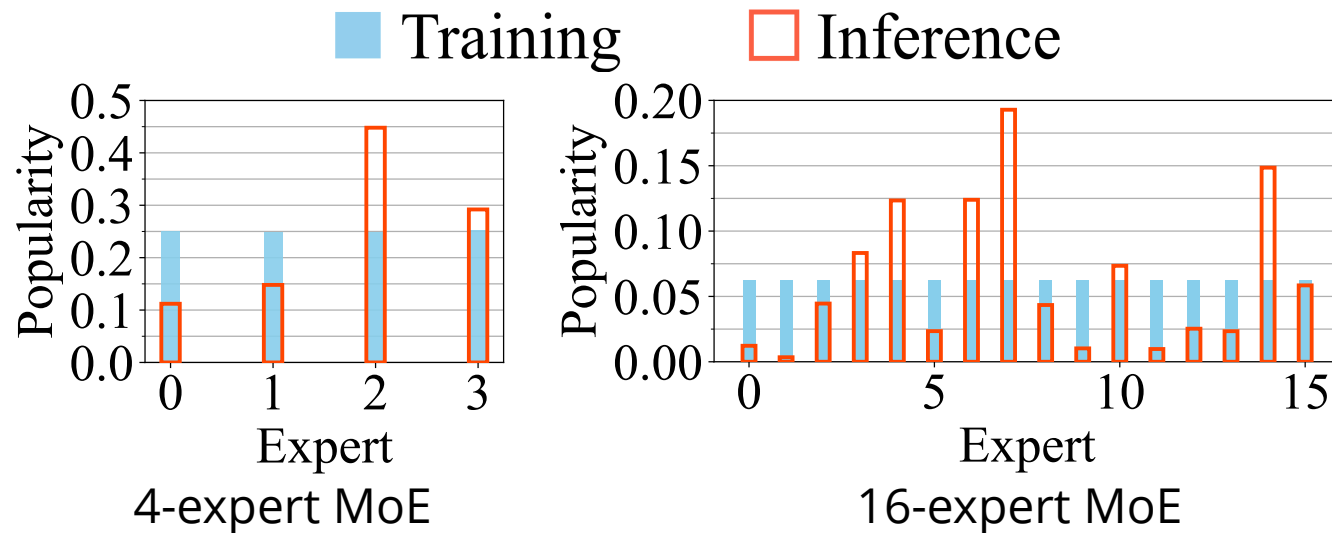
Slowdown of all-to-all varies:

Median: 2x; Maximum: ~4x



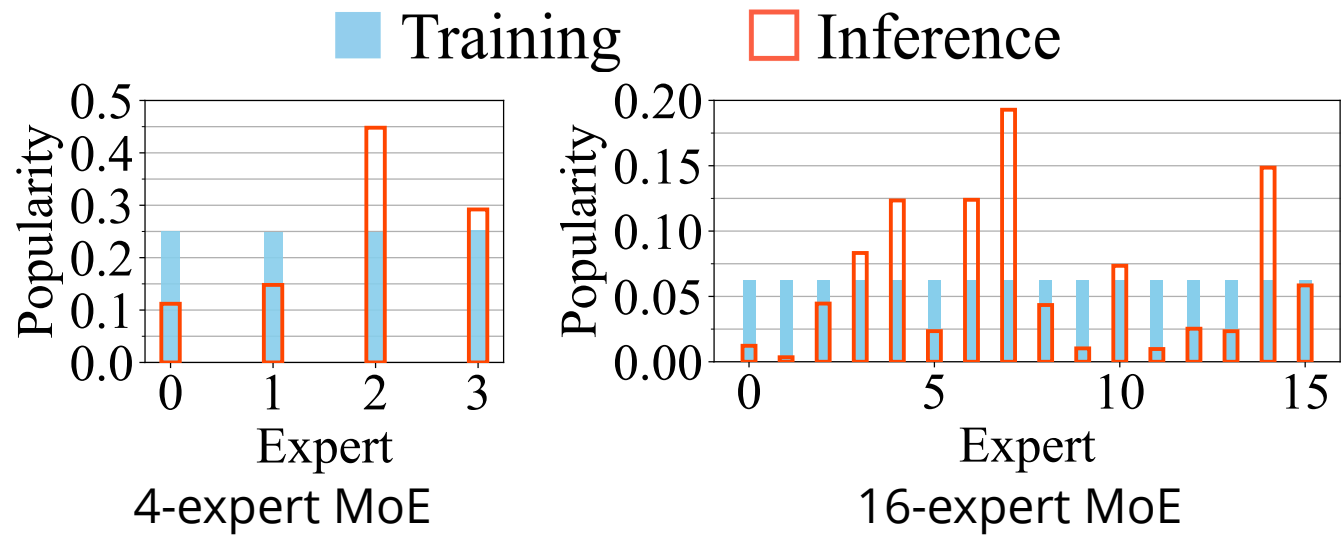
Expert Popularity in MoE Inference

- Inference: no load balancing constraints => expert selection is workload-driven, therefore, much more *biased*.



Expert Popularity in MoE Inference

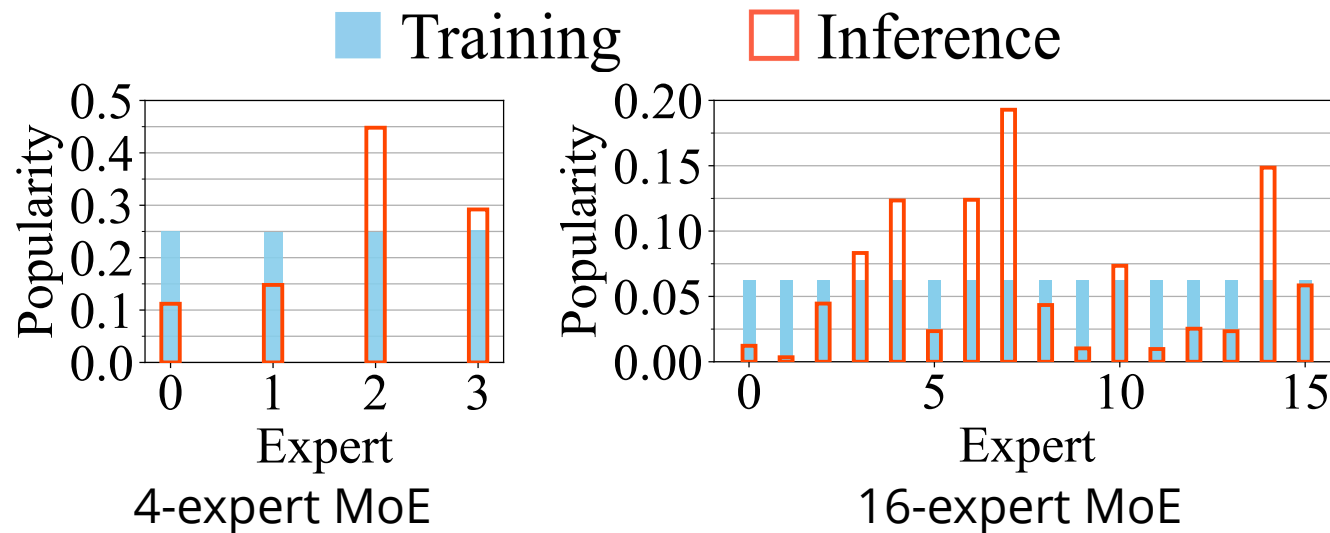
- Inference: no load balancing constraints => expert selection is workload-driven, therefore, much more *biased*.



Devices with unpopular experts have to wait for those with popular experts.

Expert Popularity in MoE Inference

- Inference: no load balancing constraints => expert selection is workload-driven, therefore, much more *biased*.



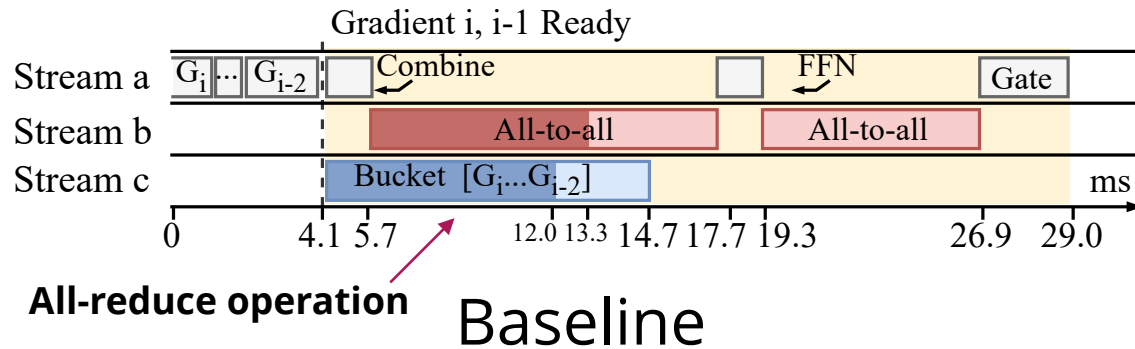
Maximum idle time of the least popular expert
=> 29.4% of the inference time.

Devices with unpopular experts have to wait for those with popular experts.

Training: Challenges

Intuition: always prioritise all-to-all and avoid bandwidth sharing.

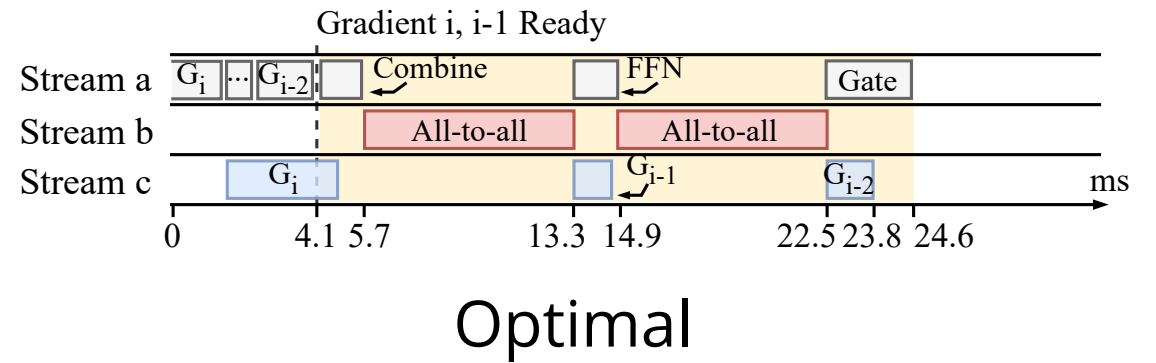
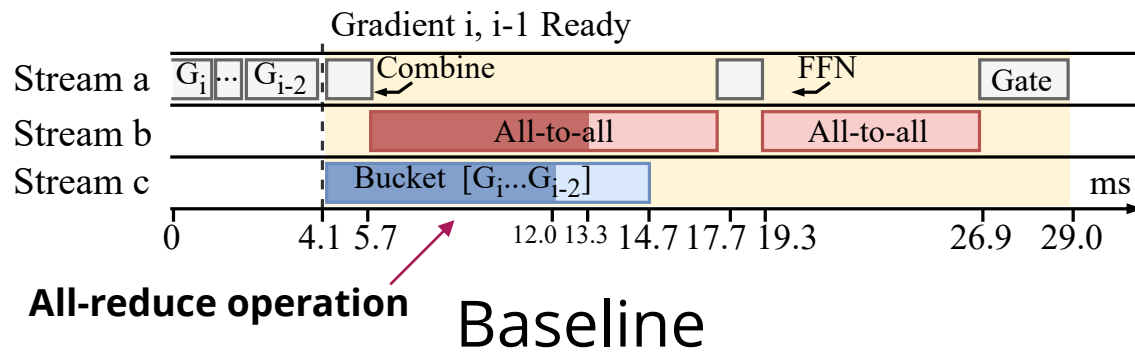
- Minimise the blocking period incurred by all-to-all



Training: Challenges

Intuition: always prioritise all-to-all and avoid bandwidth sharing.

- Minimise the blocking period incurred by all-to-all

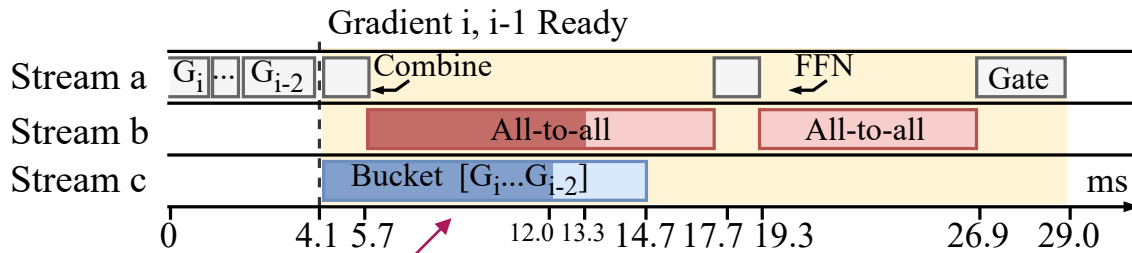


Training: Challenges

Intuition: always prioritise all-to-all and avoid bandwidth sharing.

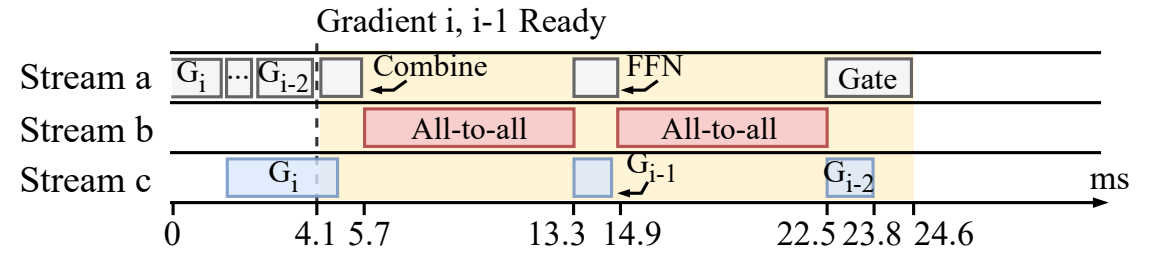
- Minimise the blocking period incurred by all-to-all

Backward pass



All-reduce operation

Baseline



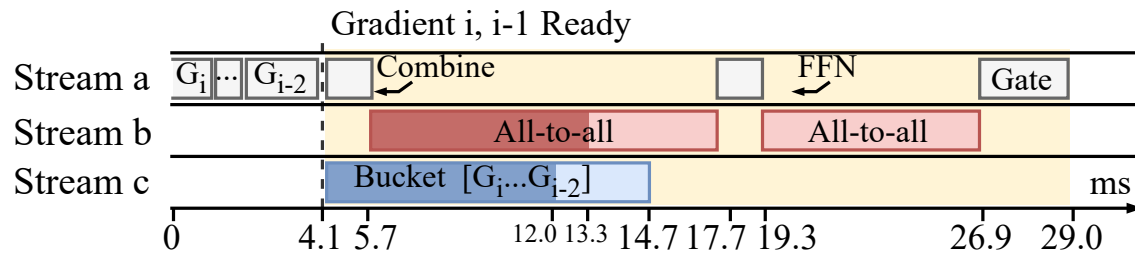
Optimal

Challenges:

1. NCCL Communication primitives cannot be preempted.
2. No control knob to adjust resource sharing (GPU SM, network bandwidth...).

Training: Micro-op Scheduling

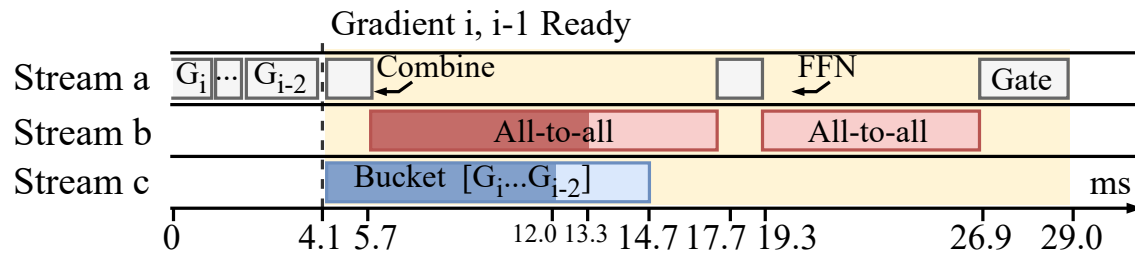
- Tensor Partitioning
 - Partition allreduce into micro-ops
 - Prioritise all-to-all whenever possible



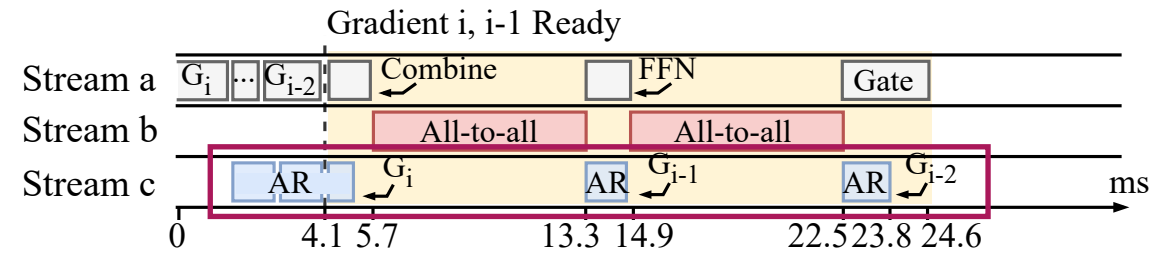
Baseline

Training: Micro-op Scheduling

- Tensor Partitioning
 - Partition allreduce into micro-ops
 - Prioritise all-to-all whenever possible



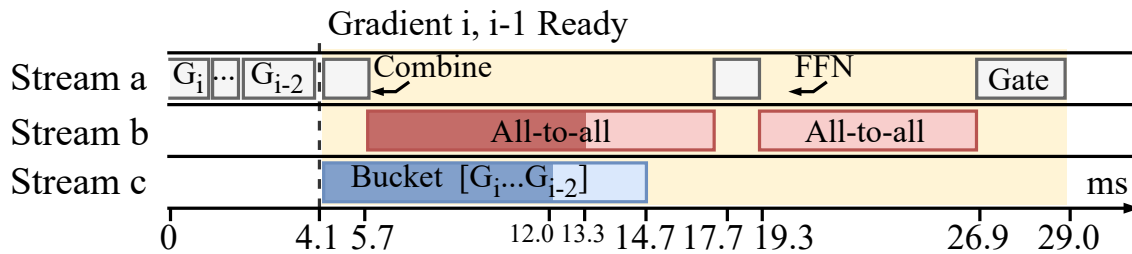
Baseline



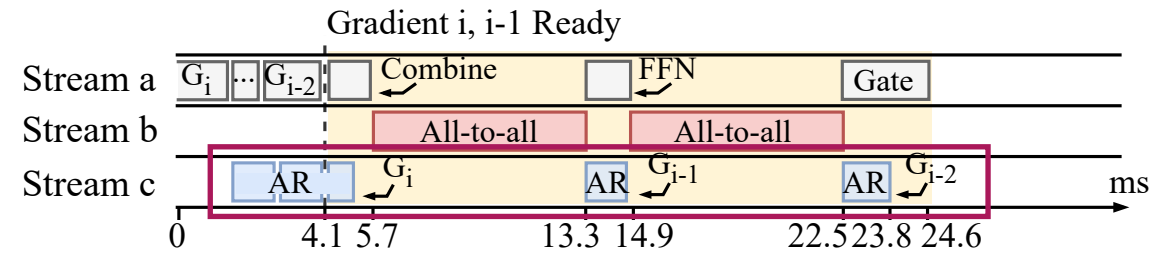
Prioritise all-to-all

Training: Micro-op Scheduling

- Tensor Partitioning
 - Partition allreduce into micro-ops
 - Prioritise all-to-all whenever possible



Baseline

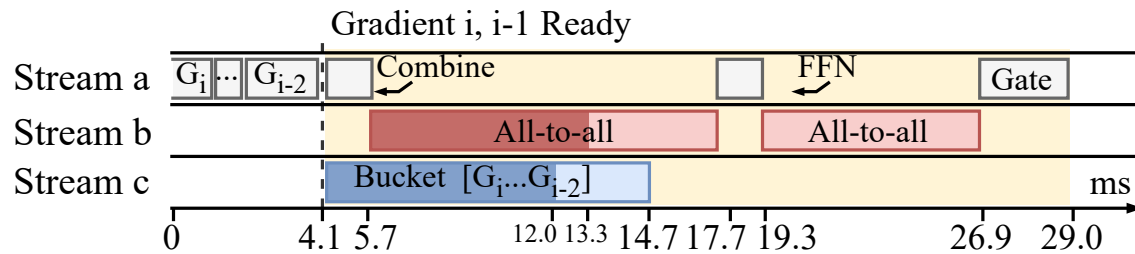


Prioritise all-to-all

- Partition all-to-all into micro-ops
- Pipelining computation and all-to-all

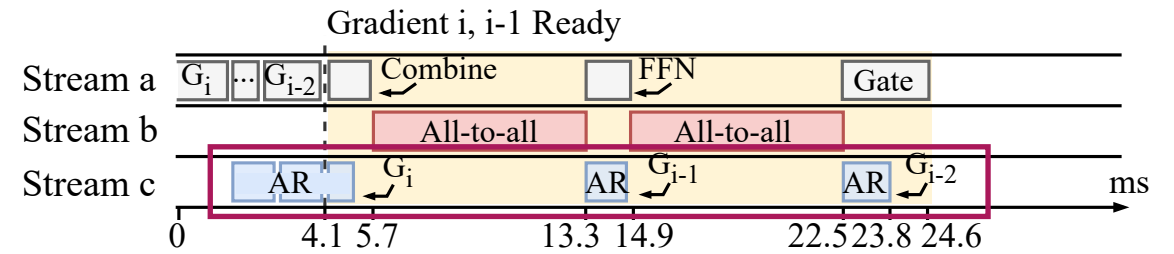
Training: Micro-op Scheduling

- Tensor Partitioning
 - Partition allreduce into micro-ops
 - Prioritise all-to-all whenever possible

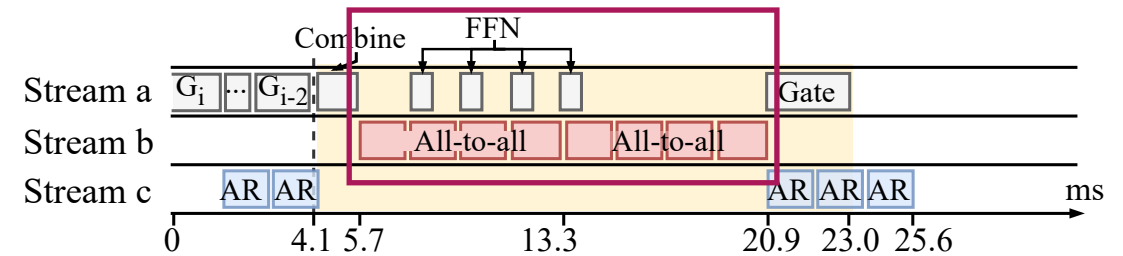


Baseline

- Partition all-to-all into micro-ops
- Pipelining computation and all-to-all

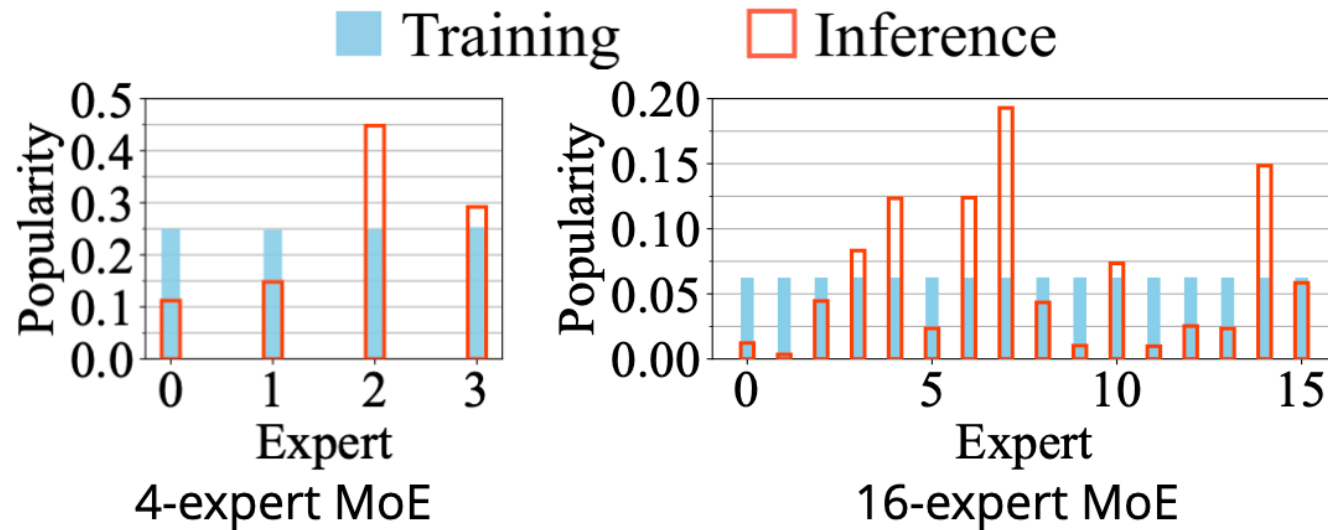


Prioritise all-to-all



Inference: Challenges

- MoE inference: no load balancing constraints => expert selection is workload-driven, therefore, much more *biased*.



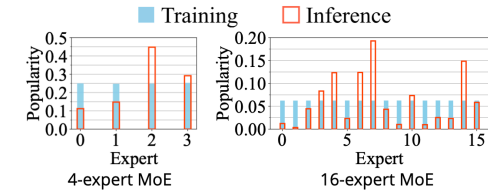
Maximum idle time of the least popular expert
=> 29.4% of the inference time.

Devices with unpopular experts have to wait for those with popular experts.

Inference: Challenges

Q: How to deal with imbalance computation load?

- MoE inference: no load balancing constraints => expert selection is workload-driven, therefore, much more *biased*.



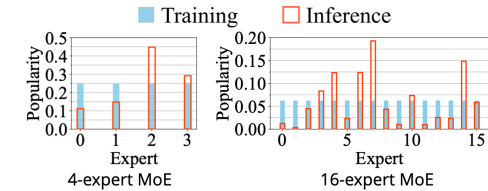
Maximum idle time of the least popular expert
=> 29.4% of the inference time.

Devices with unpopular experts have to wait for those with popular experts.

Inference: Challenges

Q: How to deal with imbalance computation load?
Allocate resources based on expert popularity:
Popular experts => more resources.

- MoE inference: no load balancing constraints => expert selection is workload-driven, therefore, much more *biased*.



Maximum idle time of the least popular expert
=> 29.4% of the inference time.

Devices with unpopular experts have to wait for those with popular experts.

Inference: Challenges

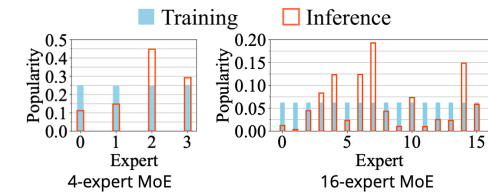
Q: How to deal with imbalance computation load?

Allocate resources based on expert popularity:

Popular experts => more resources.

Token's expert selection cannot be determined a priori to the actual gating computation.

- MoE inference: no load balancing constraints => expert selection is workload-driven, therefore, much more *biased*.



Maximum idle time of the least popular expert => 29.4% of the inference time.

Devices with unpopular experts have to wait for those with popular experts.

Model& Dataset	Layer	Top-4			
		1	2	3	4
Transformer-XL & Enwik8 (Text generation)	3	9	4	5	10
	4	5	7	8	10
	8	9	2	3	13
	12	4	5	15	8
BERT-Large & WMT En-De (Translation)	6	7	6	10	1
	8	10	6	2	15
	10	9	4	11	8
	12	1	8	10	14

Inference: Challenges

Q: How to deal with imbalance computation load?

Allocate resources based on expert popularity:

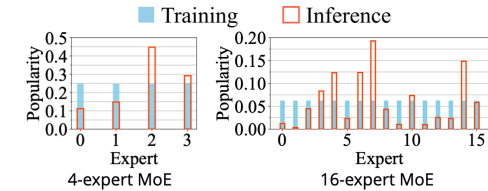
Popular experts => more resources.

Token's expert selection cannot be determined a priori to the actual gating computation.

Resource scheduling after gating network.

Inefficient practice for latency-sensitive tasks!

• MoE inference: no load balancing constraints => expert selection is workload-driven, therefore, much more *biased*.



Maximum idle time of the least popular expert => 29.4% of the inference time.

Devices with unpopular experts have to wait for those with popular experts.

Model& Dataset	Layer	Top-4			
Transformer-XL & Enwik8 (Text generation)	3	9	4	5	10
	4	5	7	8	10
	8	9	2	3	13
	12	4	5	15	8
BERT-Large & WMT En-De (Translation)	6	7	6	10	1
	8	10	6	2	15
	10	9	4	11	8
	12	1	8	10	14

Inference: Challenges

Q: How to deal with imbalance computation load?

Allocate resources based on expert popularity:

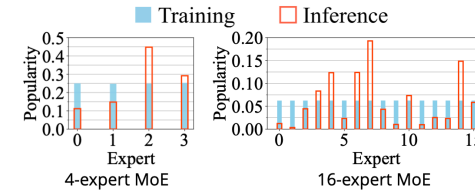
Popular experts => more resources.

Token's expert selection cannot be determined a priori to the actual gating computation.

Resource scheduling after gating network.

Inefficient practice for latency-sensitive tasks!

• MoE inference: no load balancing constraints => expert selection is workload-driven, therefore, much more *biased*.



Maximum idle time of the least popular expert => 29.4% of the inference time.

Devices with unpopular experts have to wait for those with popular experts.

Model& Dataset	Layer	Top-4			
Transformer-XL & Enwik8 (Text generation)	3	9	4	5	10
	4	5	7	8	10
	8	9	2	3	13
	12	4	5	15	8
BERT-Large & WMT En-De (Translation)	6	7	6	10	1
	8	10	6	2	15
	10	9	4	11	8
	12	1	8	10	14

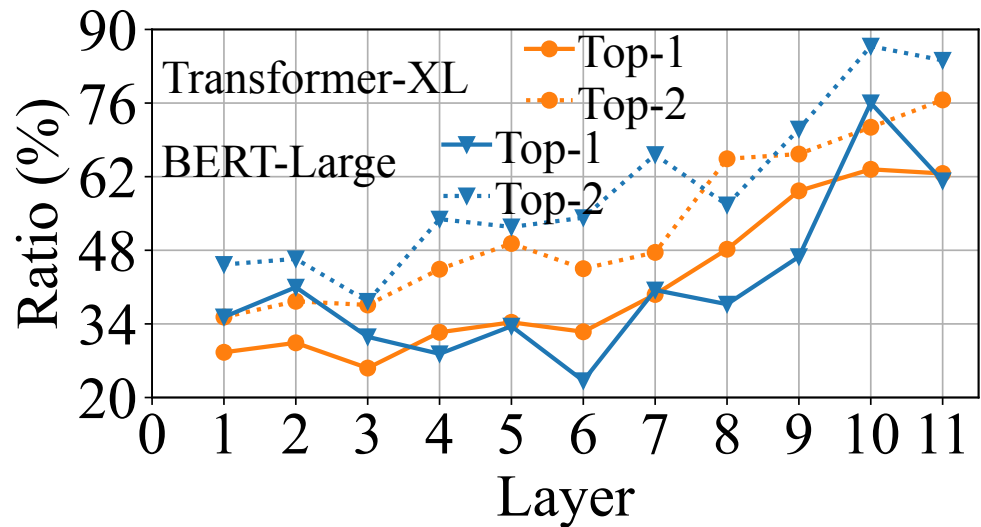
How to achieve low-overhead resource scheduling to balance device load?

Inference: Pattern in Expert Selection

- Findings: similar tokens tend to be processed by the same or similar experts in each layer.
- Tokens selecting the same expert in layer i tend to select the same expert again in layer $i + 1$.

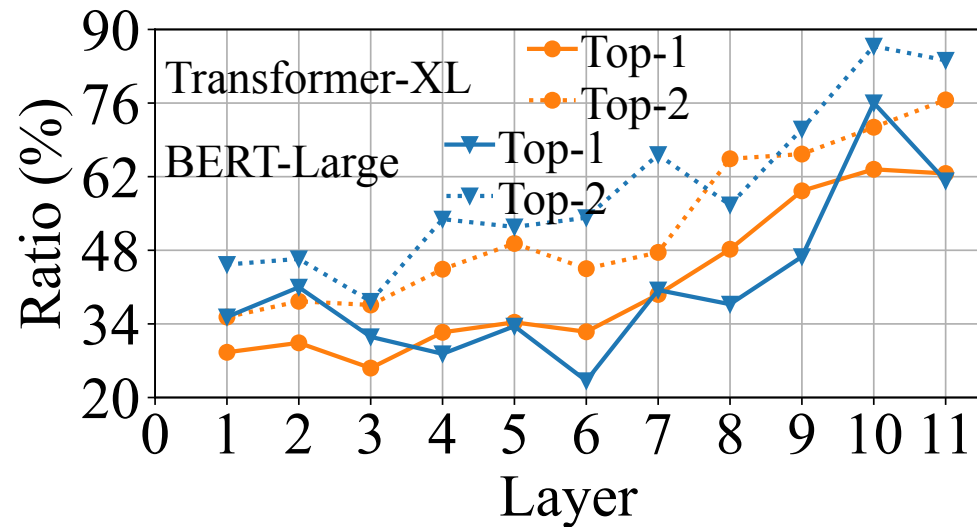
Inference: Pattern in Expert Selection

- Findings: similar tokens tend to be processed by the same or similar experts in each layer.
- Tokens selecting the same expert in layer i tend to select the same expert again in layer $i + 1$.



Inference: Pattern in Expert Selection

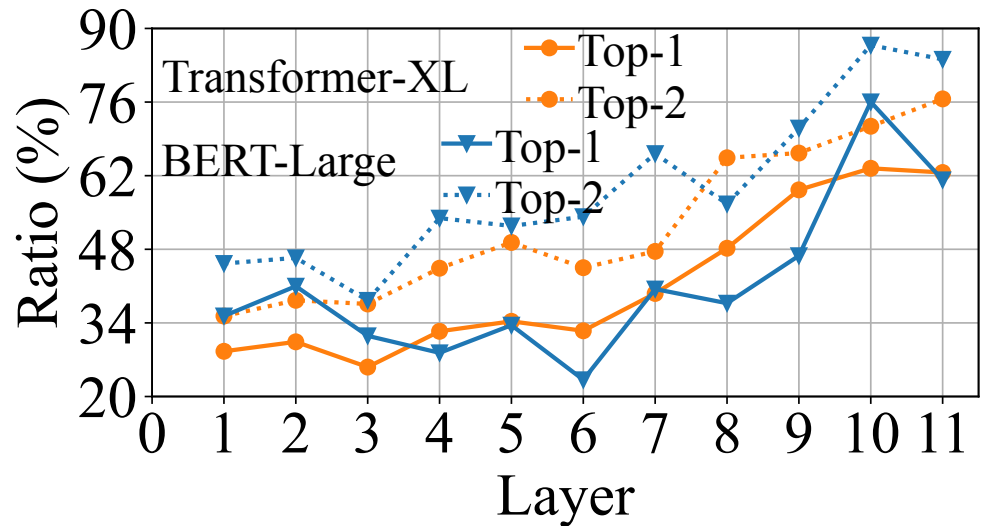
- Findings: similar tokens tend to be processed by the same or similar experts in each layer.
- Tokens selecting the same expert in layer i tend to select the same expert again in layer $i + 1$.



41.94% tokens when k is 1
54.59% tokens when k is 2

Inference: Pattern in Expert Selection

- Findings: similar tokens tend to be processed by the same or similar experts in each layer.
- Tokens selecting the same expert in layer i tend to select the same expert again in layer $i + 1$.



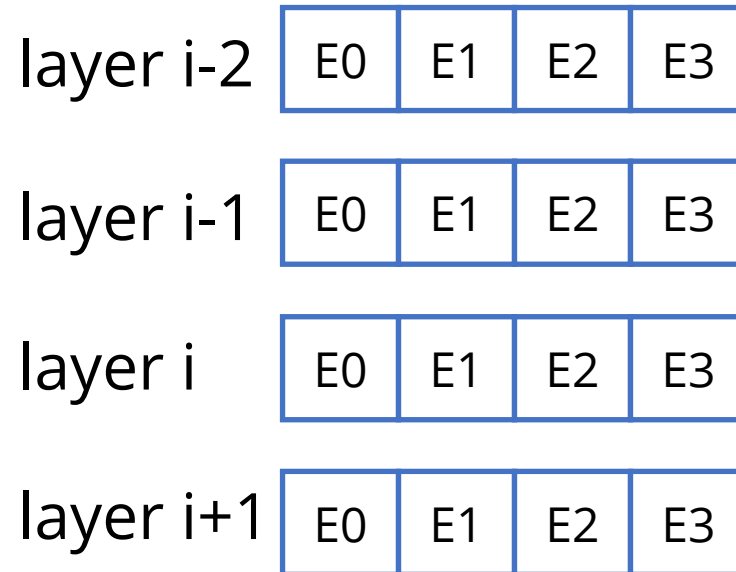
41.94% tokens when k is 1
54.59% tokens when k is 2

Exploit this pattern to estimate the overall expert popularity

Inference: Expert Popularity Estimation

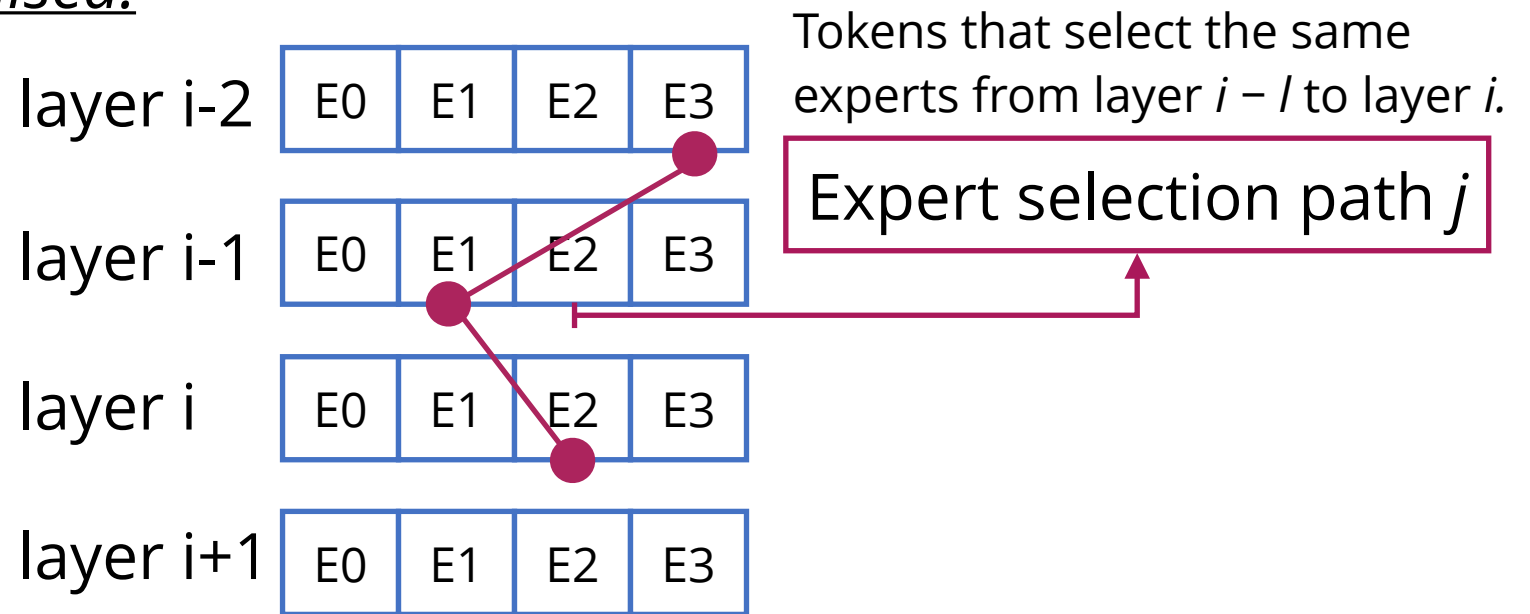
Inference: Expert Popularity Estimation

- Idea: Collect the expert selection distribution during training after the load balancing loss is minimised.



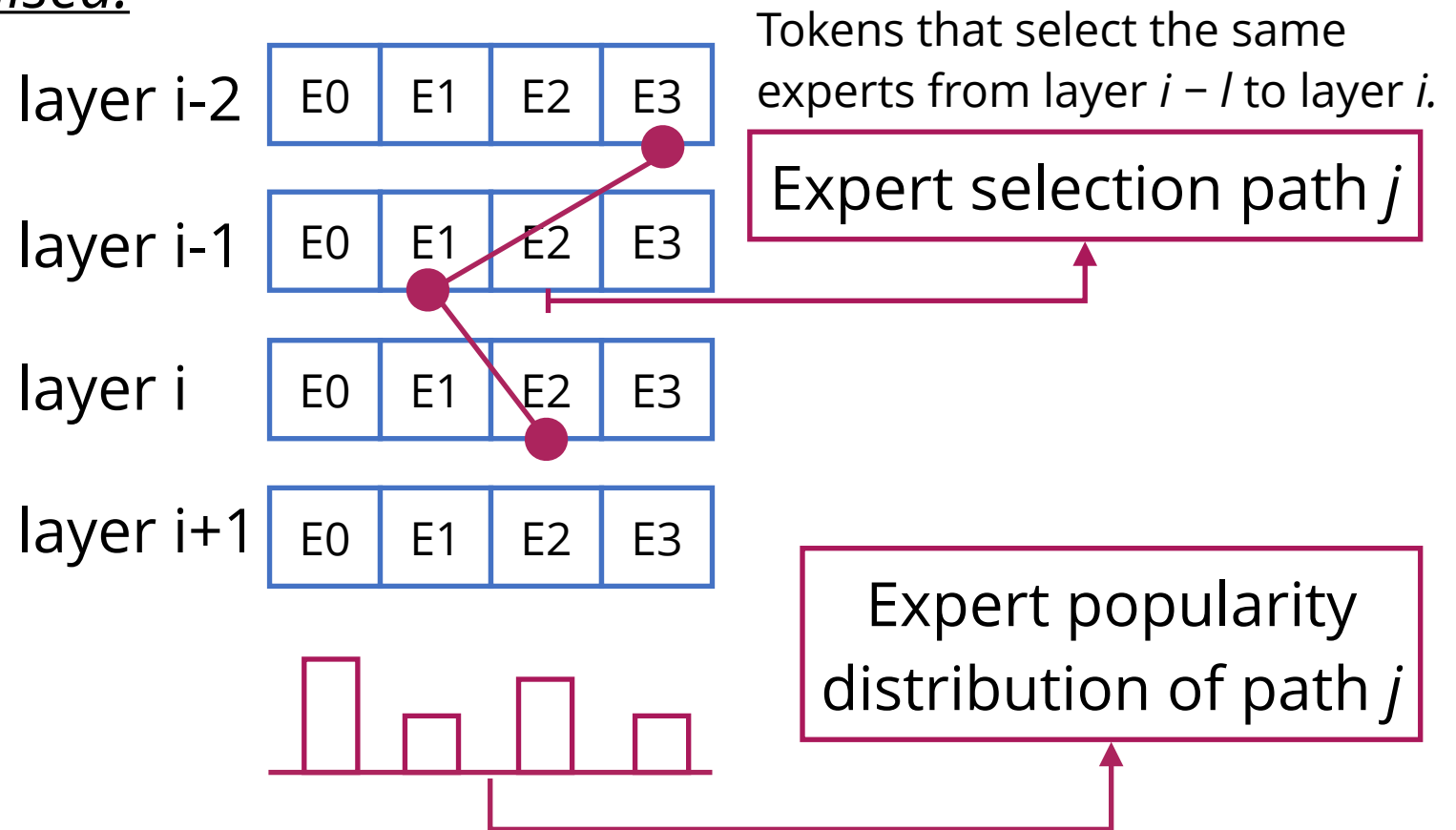
Inference: Expert Popularity Estimation

- Idea: Collect the expert selection distribution during training after the load balancing loss is minimised.



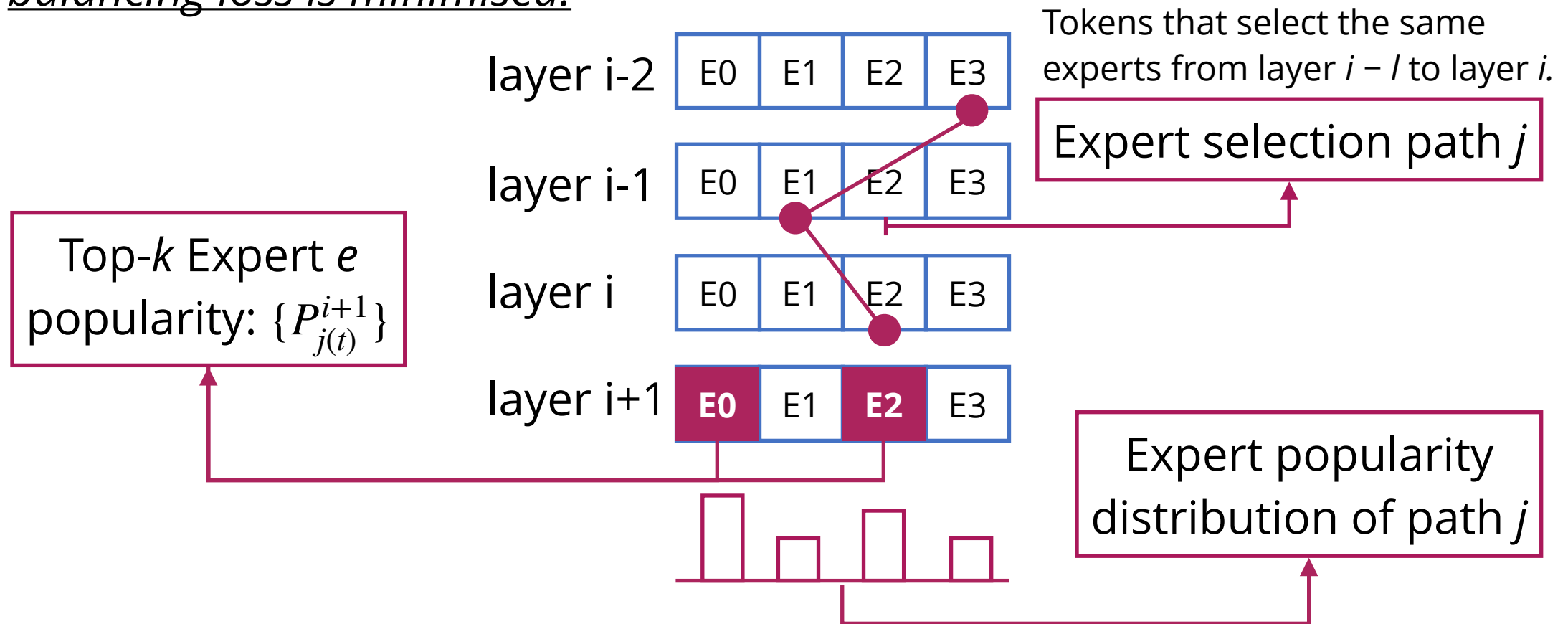
Inference: Expert Popularity Estimation

- Idea: Collect the expert selection distribution during training after the load balancing loss is minimised.



Inference: Expert Popularity Estimation

- Idea: Collect the expert selection distribution during training after the load balancing loss is minimised.



Inference: Two-phase Scheduling

- Phase 1:
 - Compute resource allocation for expert e based on popularity estimation:

$$n_e = N \times \sum_{t=1}^{N_t} P_{j(t)}^{i+1}(e) / N_t$$

No. of tokens in a batch ←

No. of GPUs ←

Inference: Two-phase Scheduling

- Phase 1: **Pipelined with model computation => nearly zero overhead**
 - Compute resource allocation for expert e based on popularity estimation:

$$n_e = N \times \sum_{t=1}^{N_t} P_{j(t)}^{i+1}(e) / N_t$$

No. of tokens in a batch ←

No. of GPUs ←

Inference: Two-phase Scheduling

- Phase 1: **Pipelined with model computation => nearly zero overhead**
 - Compute resource allocation for expert e based on popularity estimation:

$$n_e = N \times \sum_{t=1}^{N_t} P_{j^{(t)}}^{i+1}(e) / N_t$$

No. of tokens in a batch ←

No. of GPUs ←

- Phase 2:
 - Fine-tune the allocation with the actual expert selection.
 - Re-compute the allocation when the actual selection deviates significantly from the estimation.

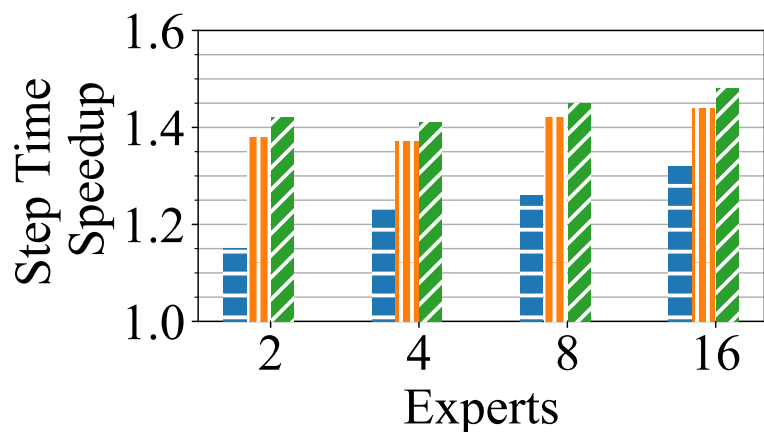
Evaluation

- Testbed: Our testbed has four worker nodes. Each node has 4 Ampere A100 GPUs with 40GB memory and is equipped with 100Gbps InfiniBand.
- Every FFN layer in Transformer is replaced with the MoE layer.
- Training models:
 - Transformer-XL: a 24-layer encoder model.
 - BERT2GPT2: a 12-layer encoder-decoder model.
 - GPT-2: a 12-layer decoder model.
- Inference models:
 - Transformer-XL: text generation with Enwik8 test set.
 - BERT: a 12-layer decoder model for translation using WMT En-De test set.

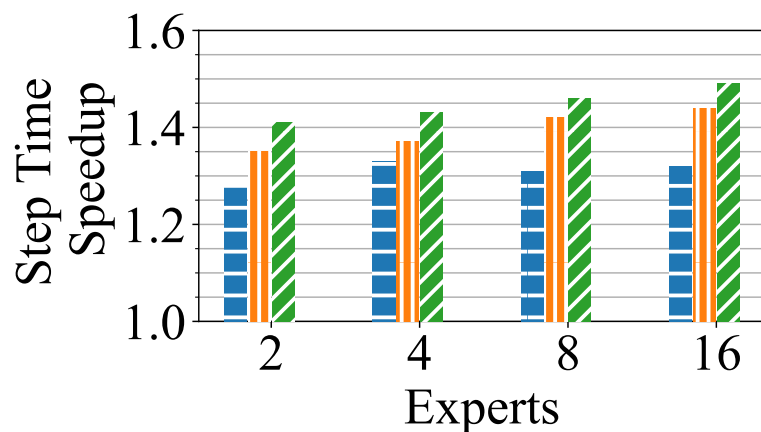
Training Step Time

Training step time speedup over Baseline (DeepSpeed) with different design choices.

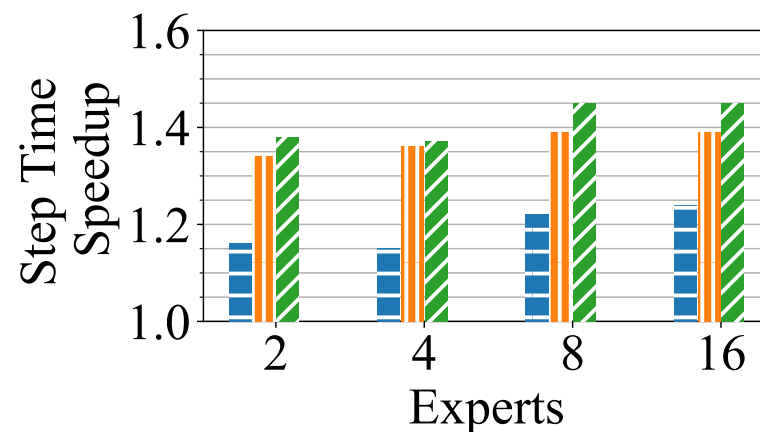
Priority only Priority + Partition Priority + Partition + Pipeline



Transformer-XL



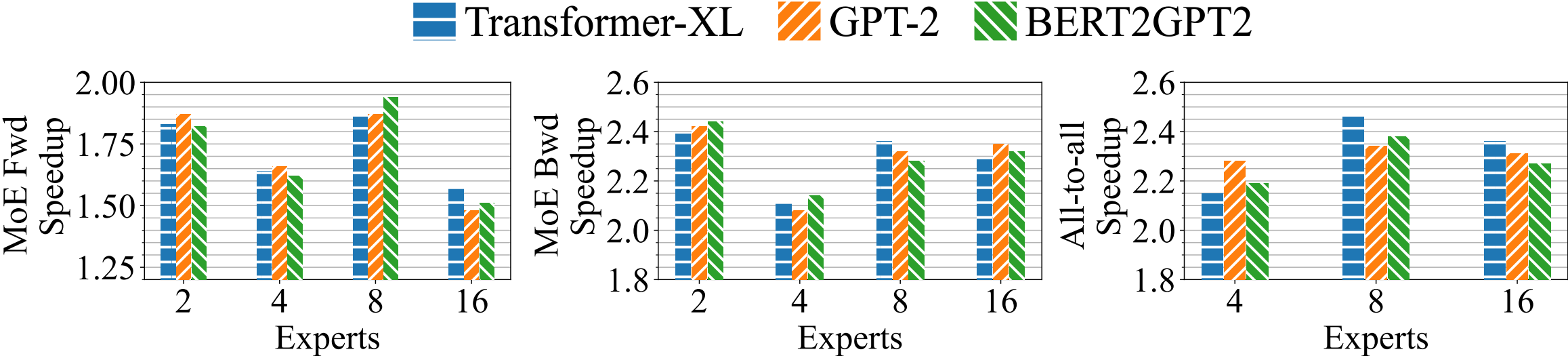
GPT-2



BERT2GPT2

MoE Layer in Training

MoE layer forward, backward, all-to-all running time speedup over Baseline.

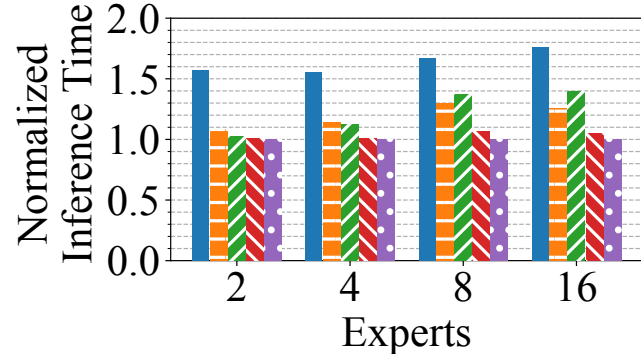
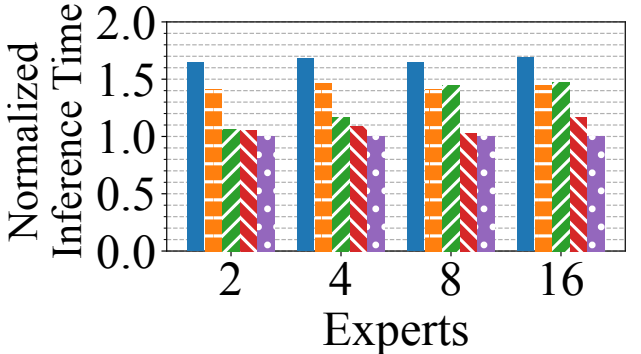


Inference Latency

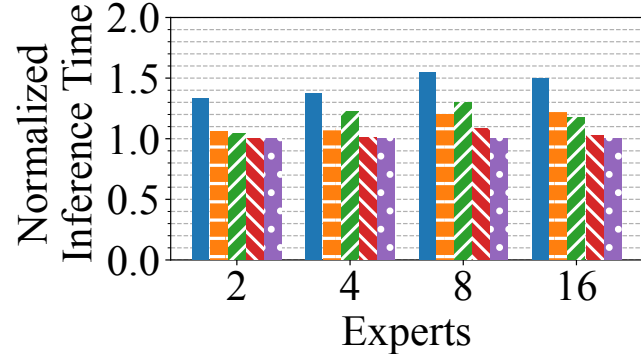
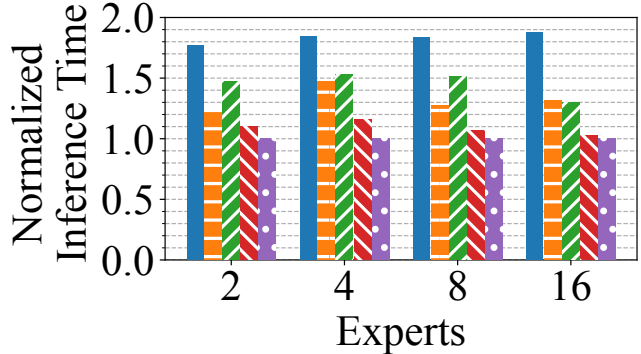
Ideal: schedule resources assuming we have the prior knowledge of exact expert popularity.

■ Baseline
 ■ Lina w/o estimation
 ■ Lina w/o fine-tuning
 ■ Lina
 ■ Ideal

Median Latency



95%ile Latency

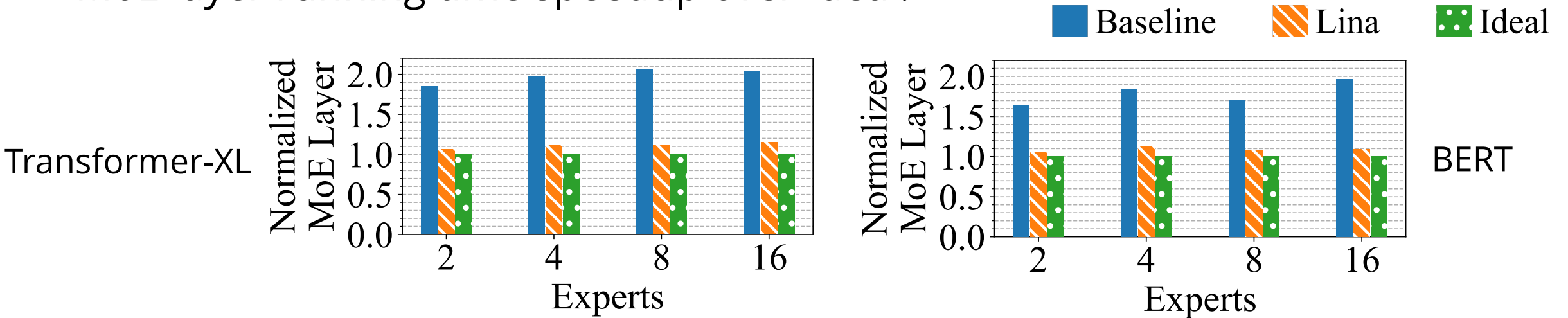


Transformer-XL

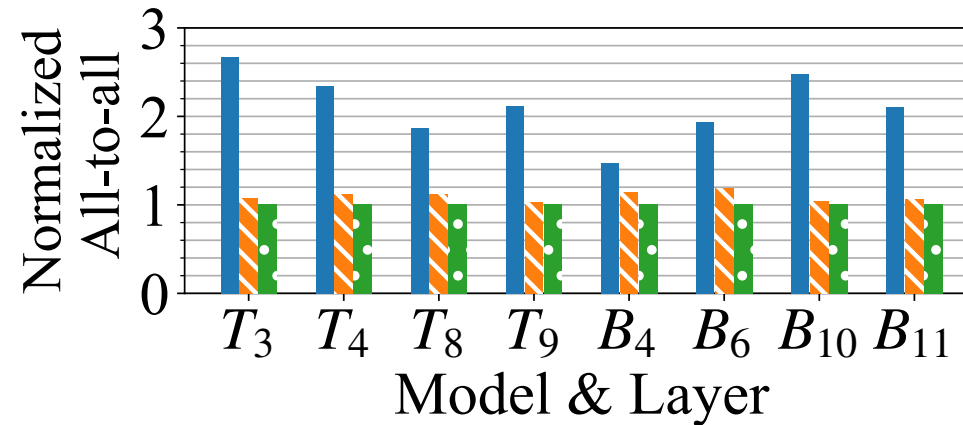
BERT

MoE Layer in Inference

MoE layer running time speedup over Ideal.



All-to-all running time speedup over Ideal in selected layers.



Lina's Contributions

- An in-depth empirical analysis of distributed MoE
 - Main causes for all-to-all to be the performance bottleneck in training and inference.
- [Training] A scheduler prioritises all-to-all over allreduce to improve its bandwidth and reduce its blocking period.
- [Inference] An estimation method of expert popularity to conduct two-phase resource scheduling.

Thanks!