



A Hybrid Tensor-Expert-Data Parallelism Approach to Optimize Mixture-of-Experts Training

Siddharth Singh¹, Olatunji Ruwase², Ammar Ahmad Awan², Samyam Rajbhandari²,

Yuxiong He², Abhinav Bhatele¹

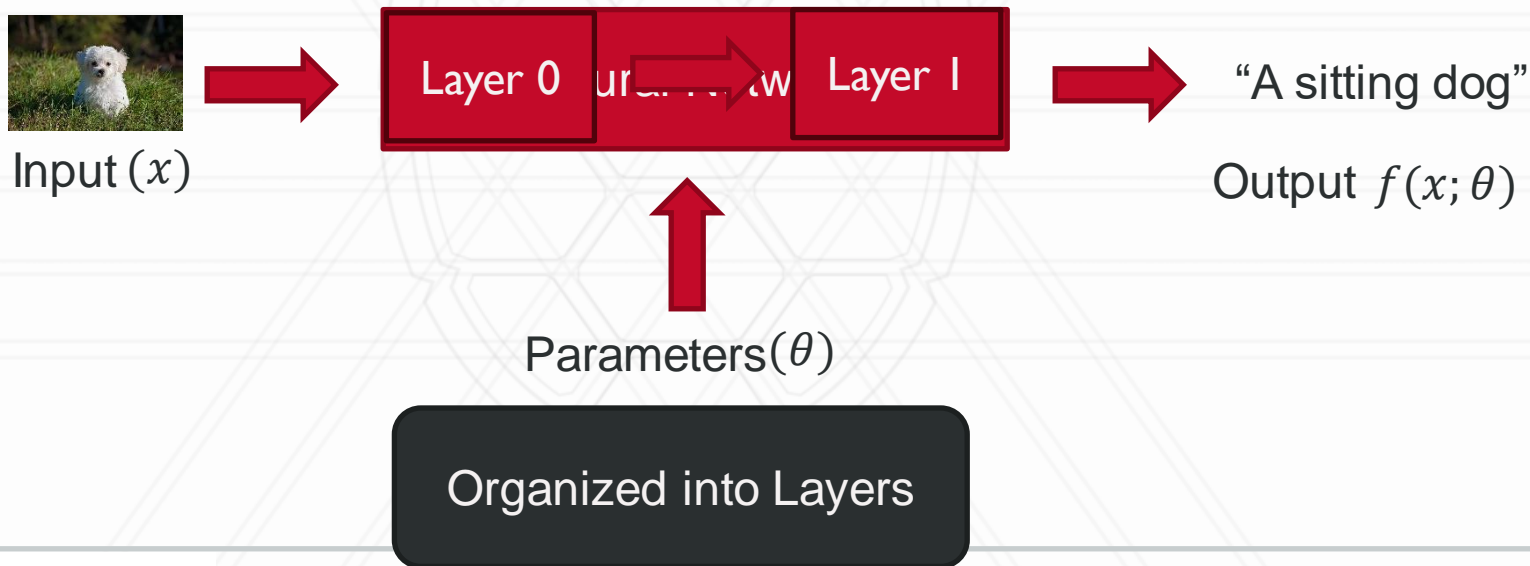
University of Maryland¹, Microsoft, Inc.²



UNIVERSITY OF
MARYLAND

Neural Networks

- Neural Networks (NNs): ‘Parameterized’ function approximators
- Can work with very high dimensional data.



Stochastic Gradient Descent

- Repeat until loss (L) has been minimized sufficiently:
 - Read in a batch of training data
 - Forward Pass : Calculate output $f(x; \theta)$ and the loss (L) on the batch.
 - Backward Pass : Calculate gradients of the loss wrt the parameters $(\frac{\partial L}{\partial \theta})$.
 - Optimizer Step : Use $\frac{\partial L}{\partial \theta}$ to update θ .

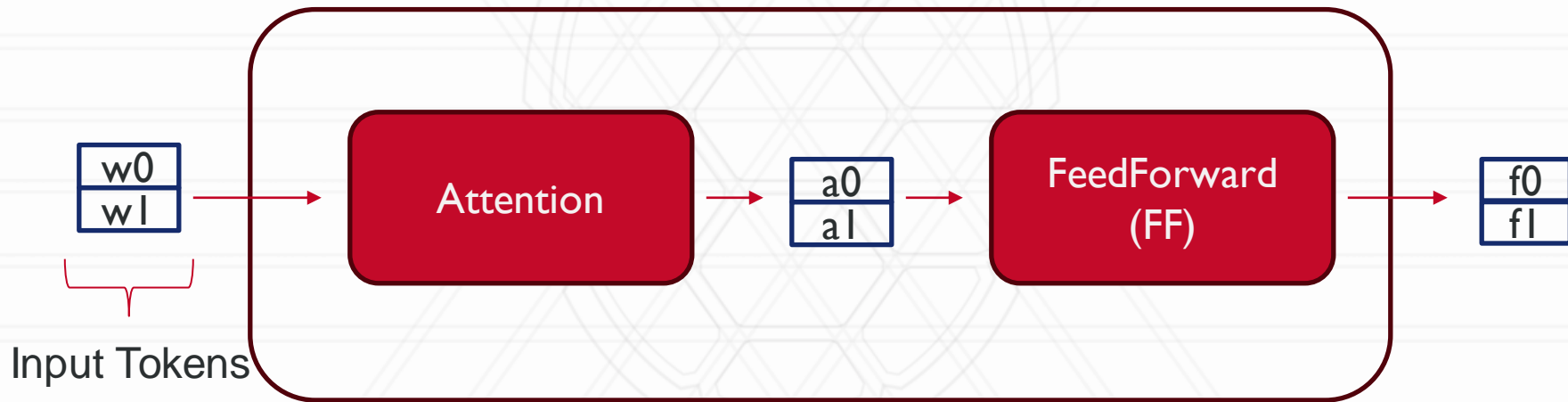
Motivation

- More parameters = More accurate neural networks 😊
- However, more parameters = More FLOPs in training 😞

MoEs can make a given model arbitrarily large without changing its training FLOPs!

Mixture-of-Experts (MoEs)

Step 1: Start with a base model – usually a transformer neural network.



A Transformer Layer (Base Model)

Mixture-of-Experts (MoEs)

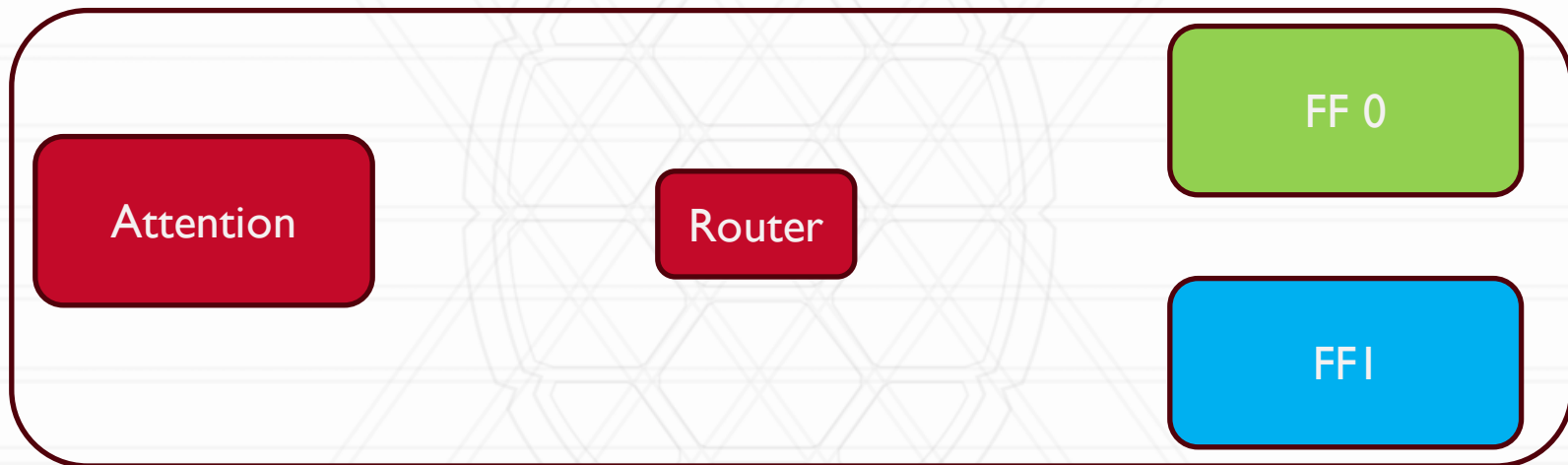
Step 2: Introduce multiple FF blocks. Each FF block is an 'expert'.



A Transformer Layer (Base Model) + 2 experts

Mixture-of-Experts (MoEs)

Step 3: Introduce a parameterized router that maps tokens to experts.

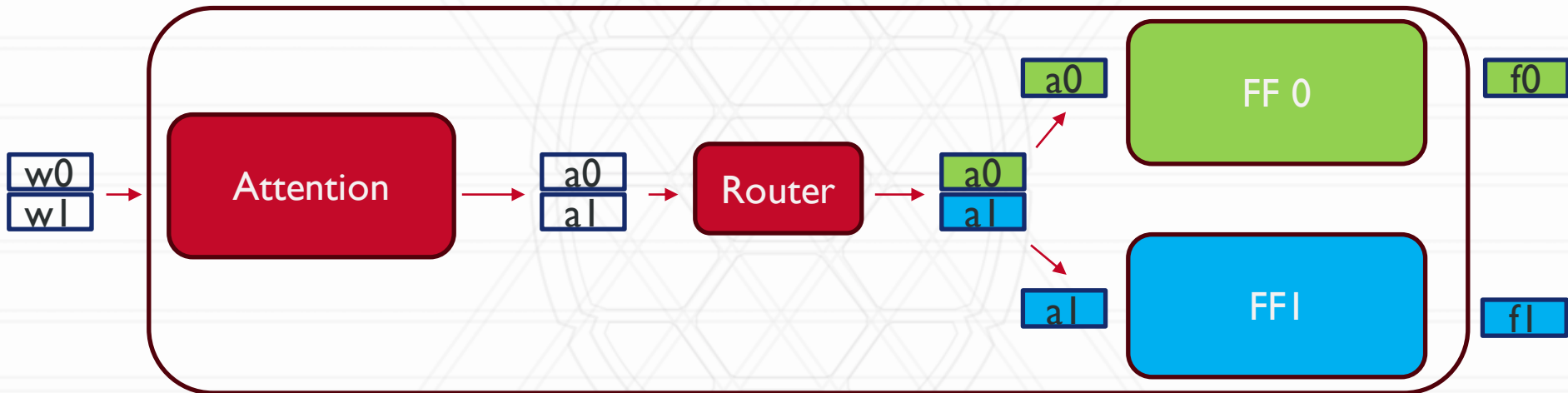


A Transformer Layer (Base Model) + 2 experts

Mixtur

Total FLOPs per token is independent of the number of experts!

Step 3: Introduce a parameterized router that maps tokens to experts.



A Transformer Layer (Base Model) + 2 experts

Caveat

- Diminishing returns beyond 64-128 experts.
- Imperative to increase base model sizes along with expert counts.

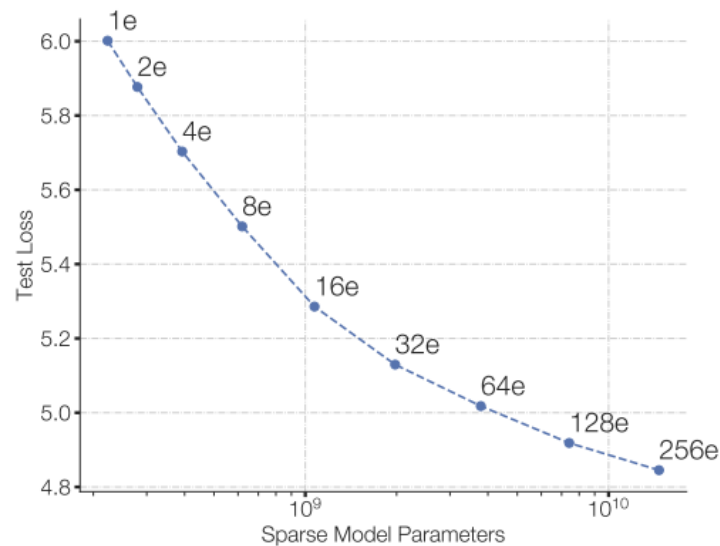


Figure courtesy Fedus et al. [2]

Gaps in Current Work

- Large MoEs are almost always trained in parallel on multiple GPUs.
- However, current parallel frameworks are not suited for MoEs built with large base models.
- They either support base models of limited sizes due to limited dimensions of parallelism.
 - Example – DeepSpeed-MoE [4] which has expert+data parallelism but not model parallelism (tensor/pipeline)
- Or, they use extremely inefficient parallel techniques like out-of-core training or FSDP.
 - Inefficiency occurs due to high communication times

Our work – Deepspeed-TED

- Goal 1 - Support MoEs with large base models
- Goal 2 - Minimize communication times to maintain efficiency.
- A three-dimensional hybrid of state-of-the-art parallel training algorithms
 - T – Tensor Parallelism (Megatron-LM [3])
 - E – Expert Parallelism (DeepSpeed-MoE [4])
 - D – Sharded Data Parallelism (ZeRO [5])

Data Parallelism

Average Gradients

GPU - 0



w0
w1

Attention

a0
a1

Router

a0
a1

a0

FF0

a0

a1

FF1

f1

w0
w1
w2
w3



w2
w3

Attention

a2
a3

Router

a2
a3

a2

FF0

f2

a3

FF1

f3

GPU - 1

Data + Expert Parallelism

Divide Experts among GPUs

GPU - 0



w0
w1

Attention

Router

a0
a1

a0
a2

FF0

FF0

FF1

f0
f2

f0
f1

w0
w1
w2
w3



w2
w3

Attention

Router

a2
a3

a1
a3

FF0

FF1

FF1

f1
f3

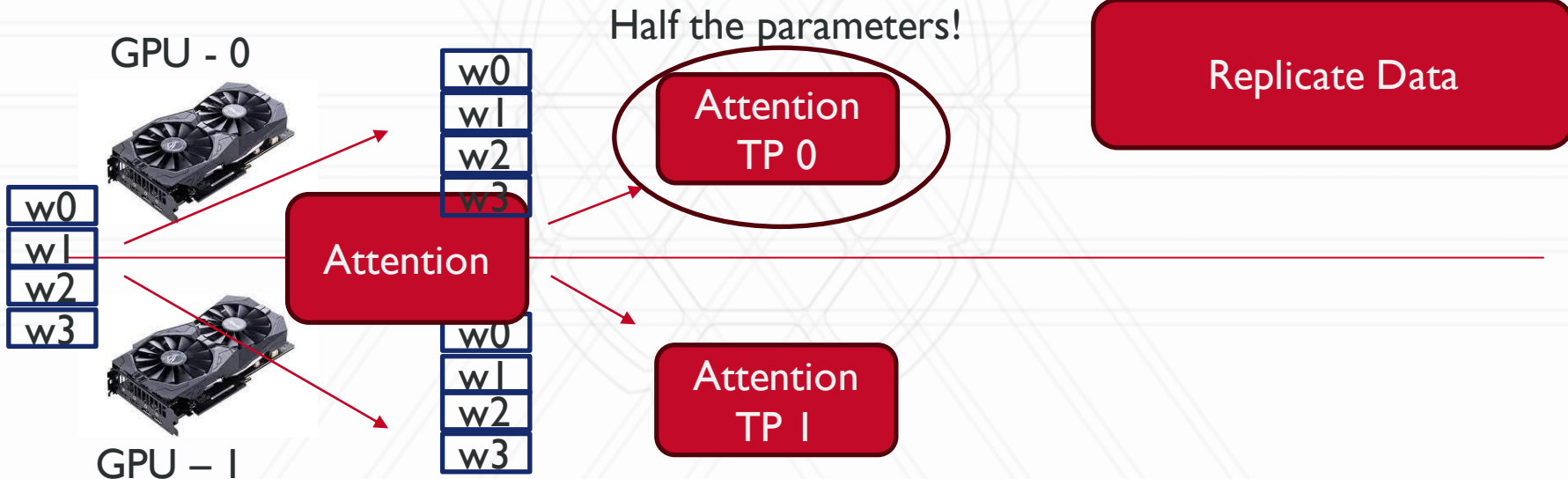
f2
f3

GPU - 1

All-to-All Communication

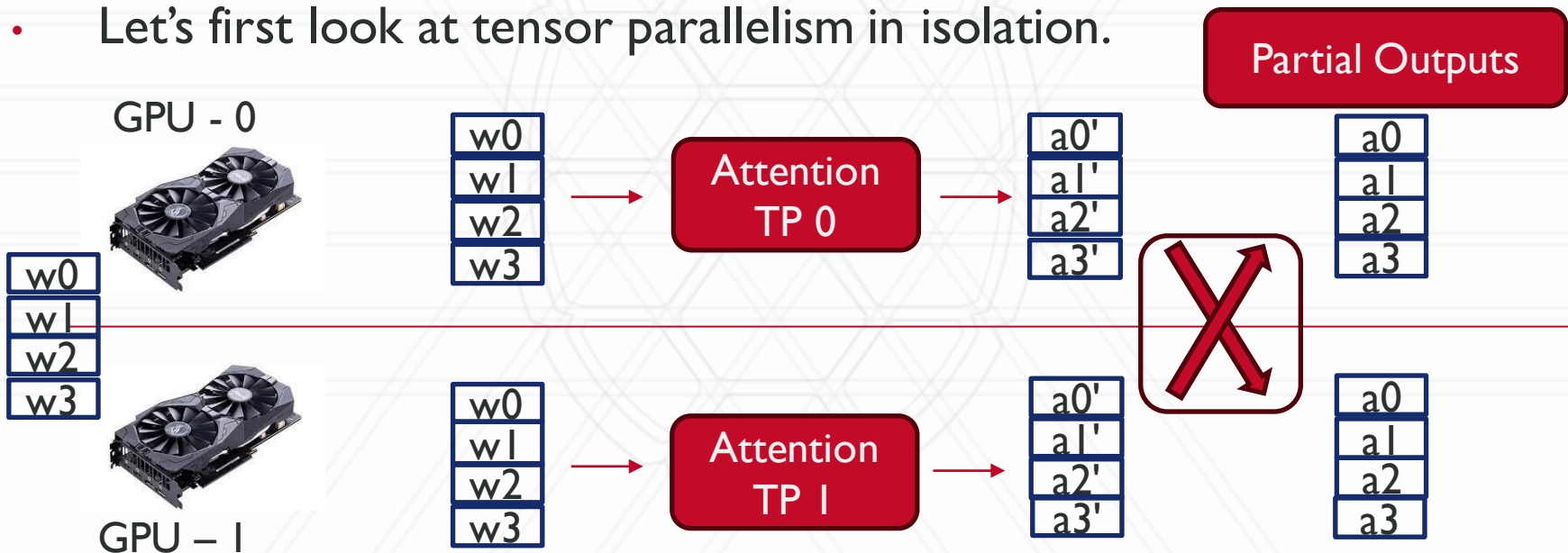
Tensor Parallelism

- Parallelize the matrix multiplications inside Attention and FF.
- Let's first look at tensor parallelism in isolation.



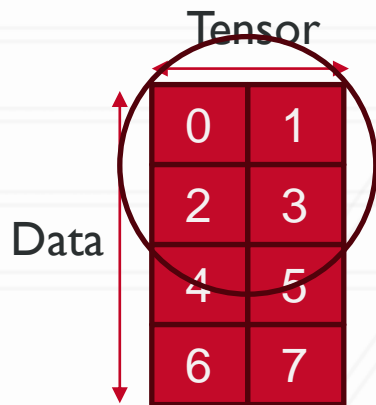
Tensor Parallelism

- Parallelize the matrix multiplications inside Attention and FF.
- Let's first look at tensor parallelism in isolation.

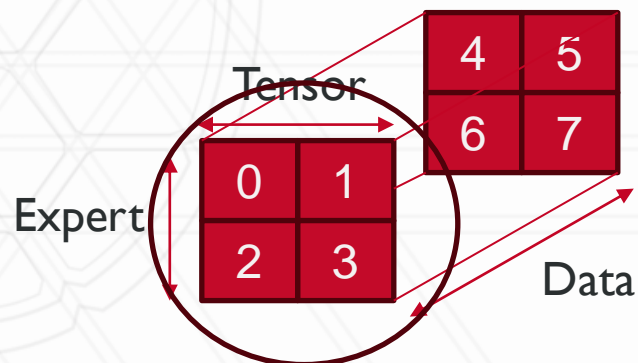


Data + Expert + Tensor Parallelism

- Now let us look at tensor+expert+data parallelism on 8 GPUs.
- Two virtual topologies for Attention and FF-Blocks.
- For brevity, we will only look at GPUs [0-3]

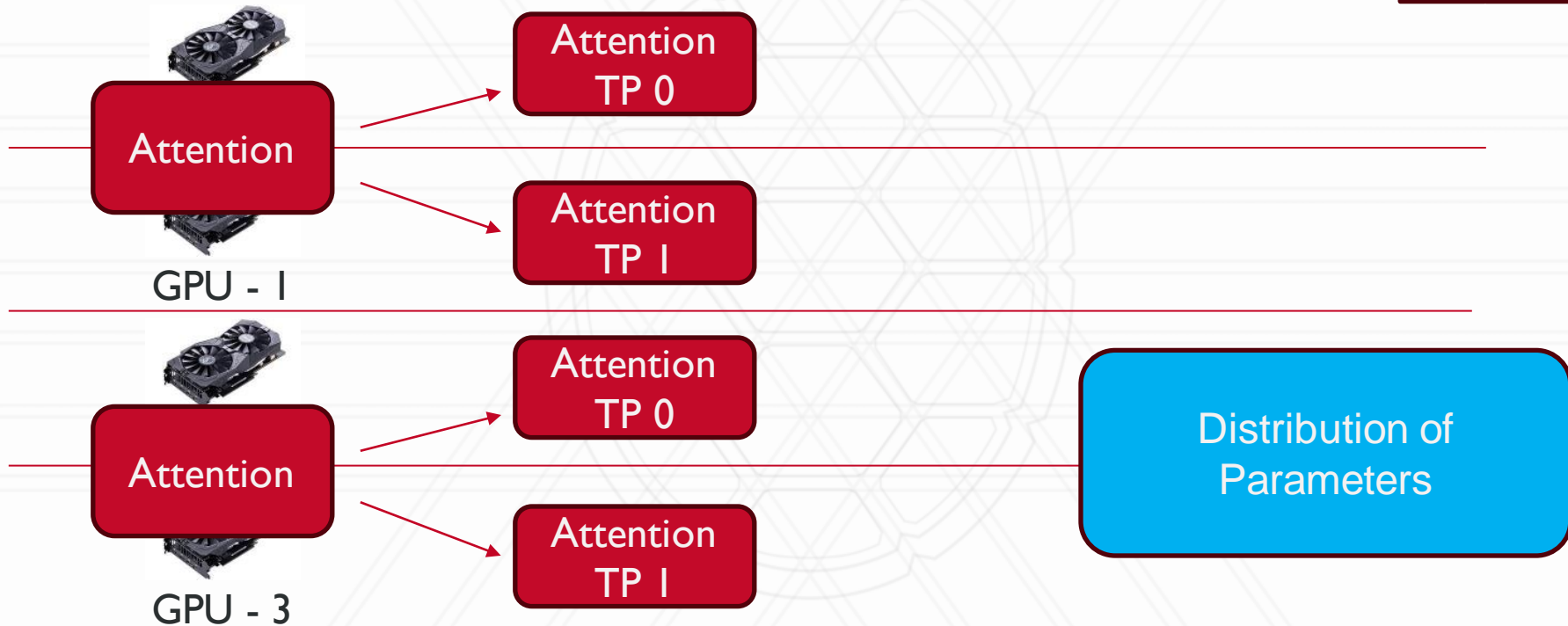
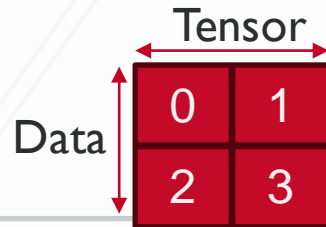


Attention Block (Non-Expert)

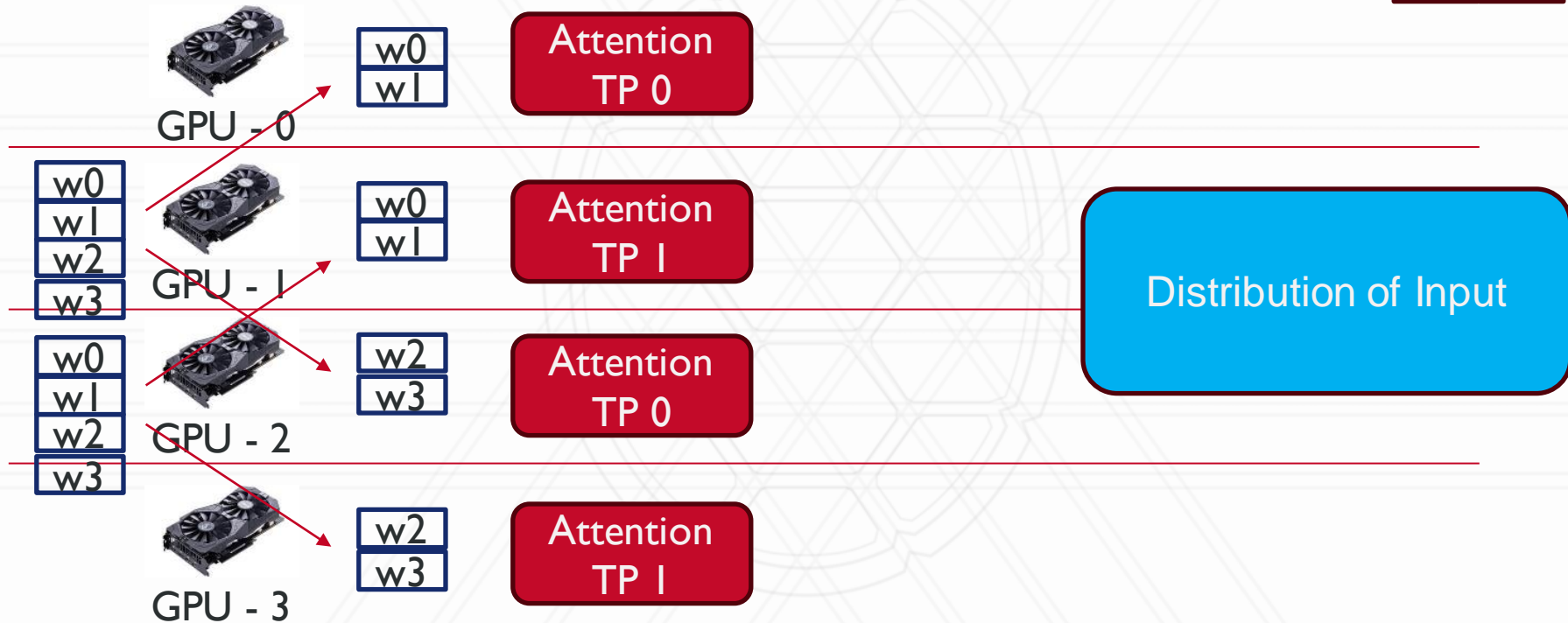
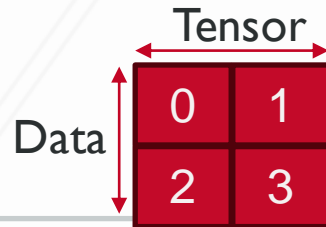


FF Blocks (Expert)

TED for Attention (Non-Experts)



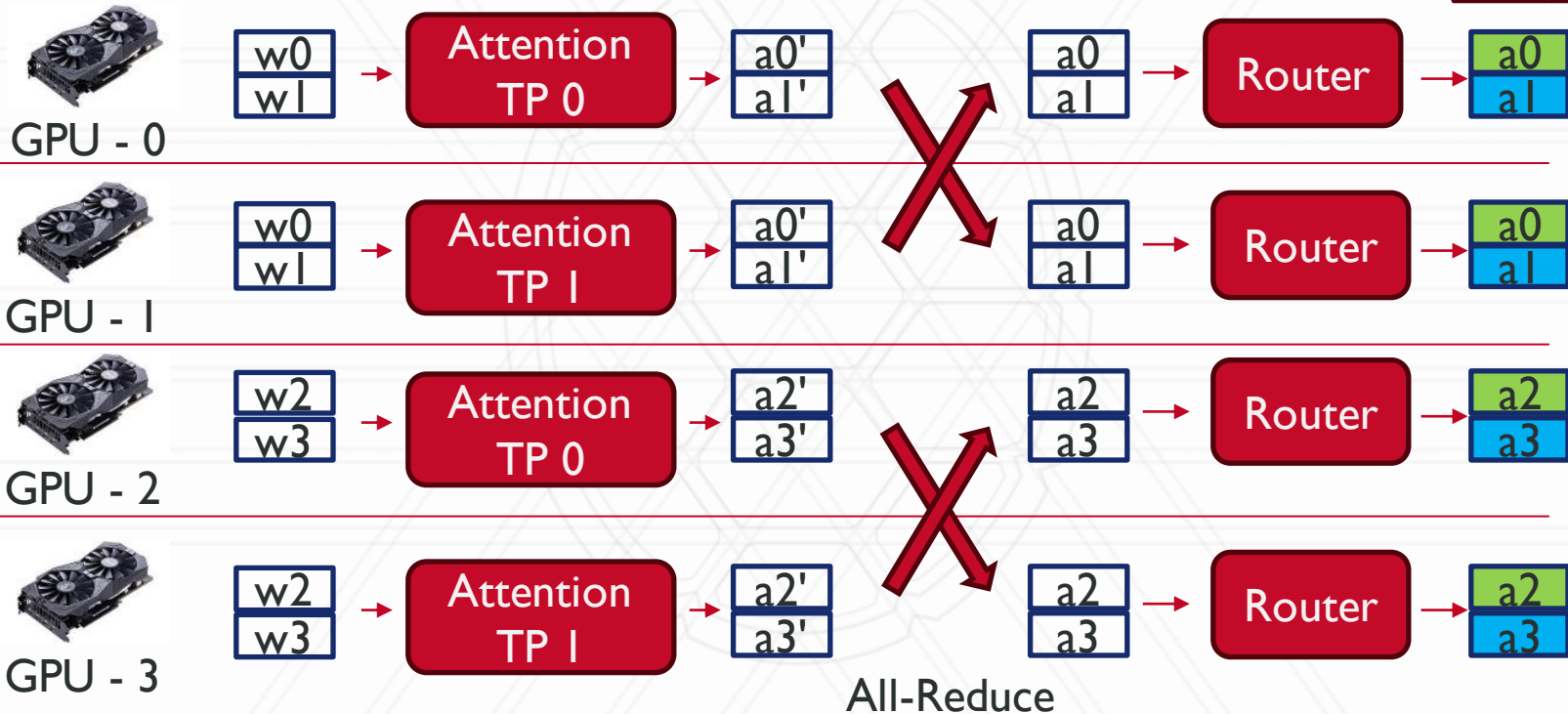
TED for Attention (Non-Experts)



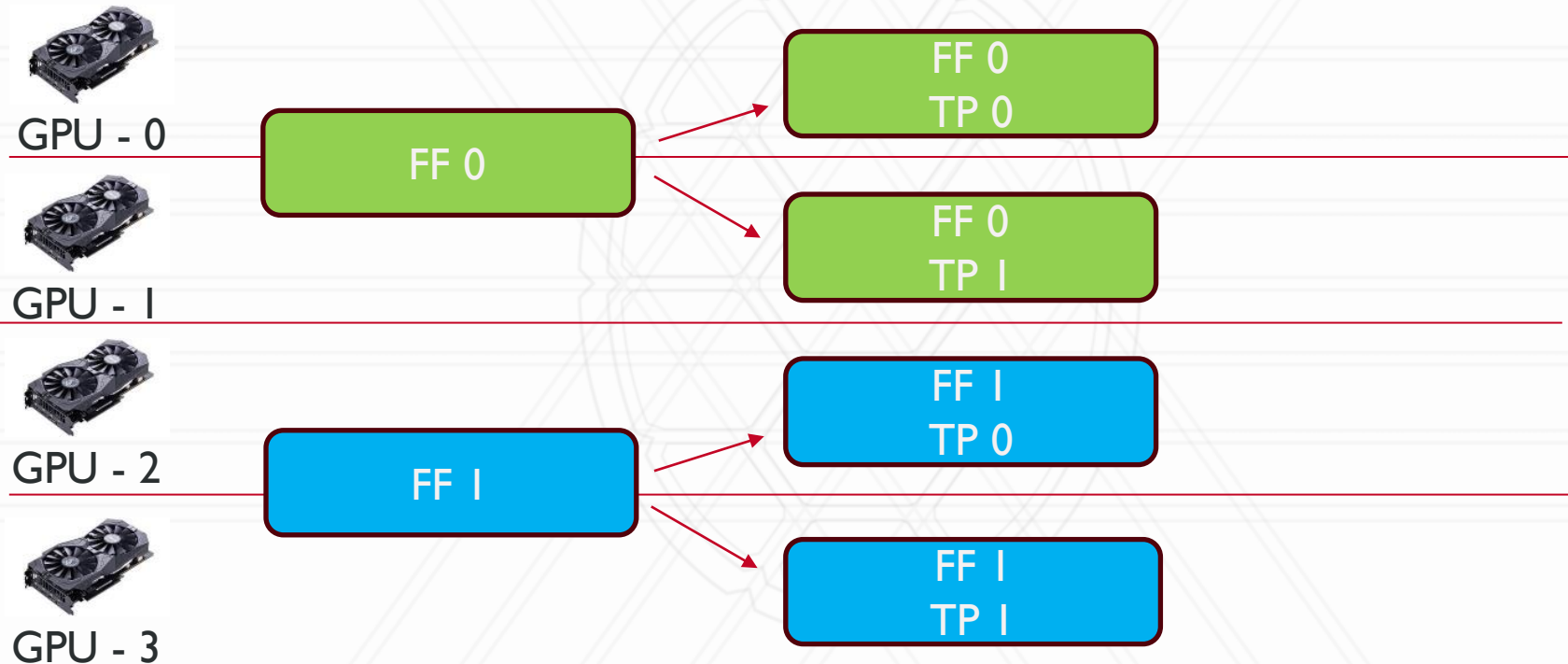
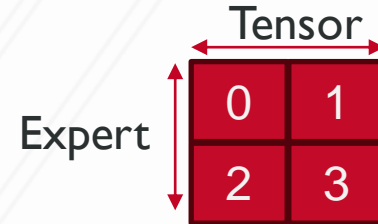
TED for Attention (Non-Experts)

Data

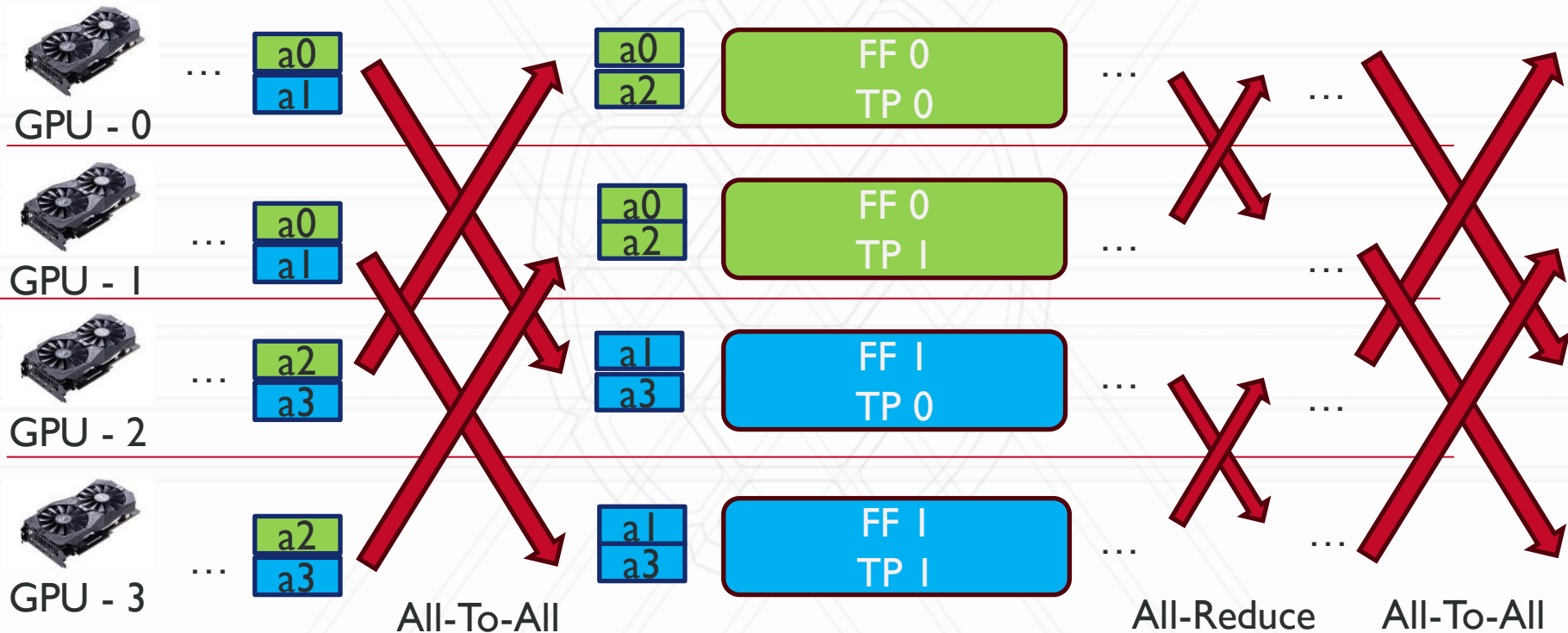
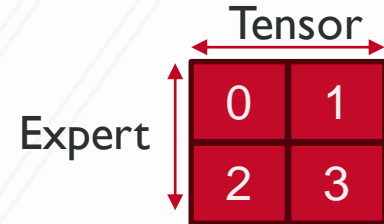
Tensor	
0	1
2	3



TED for FF (Experts)

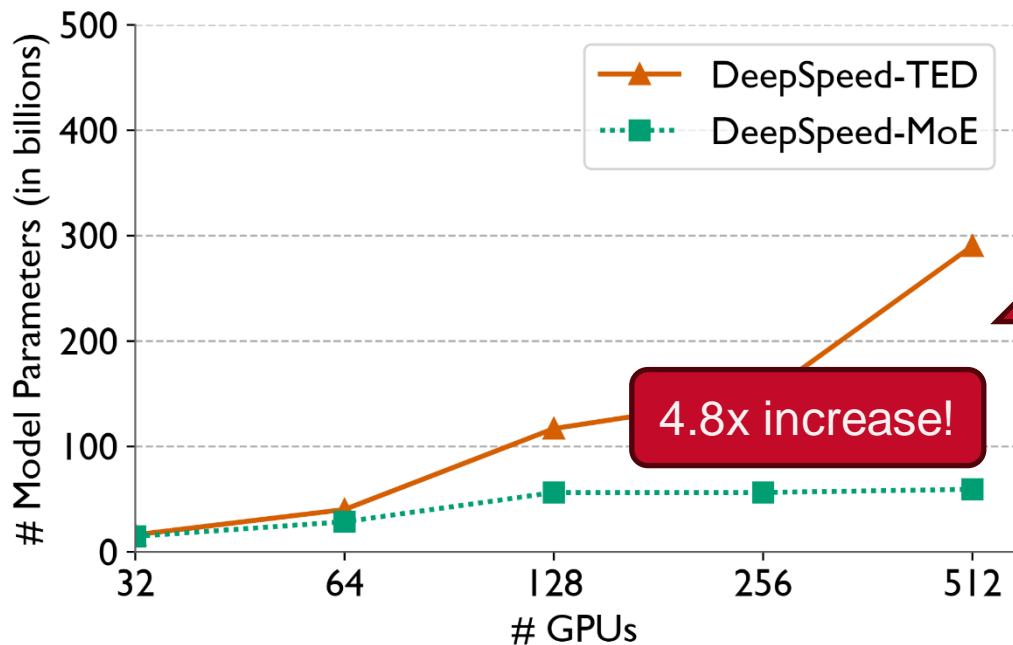


+ TED for FF (Experts)



3D parallelism helps up train larger models

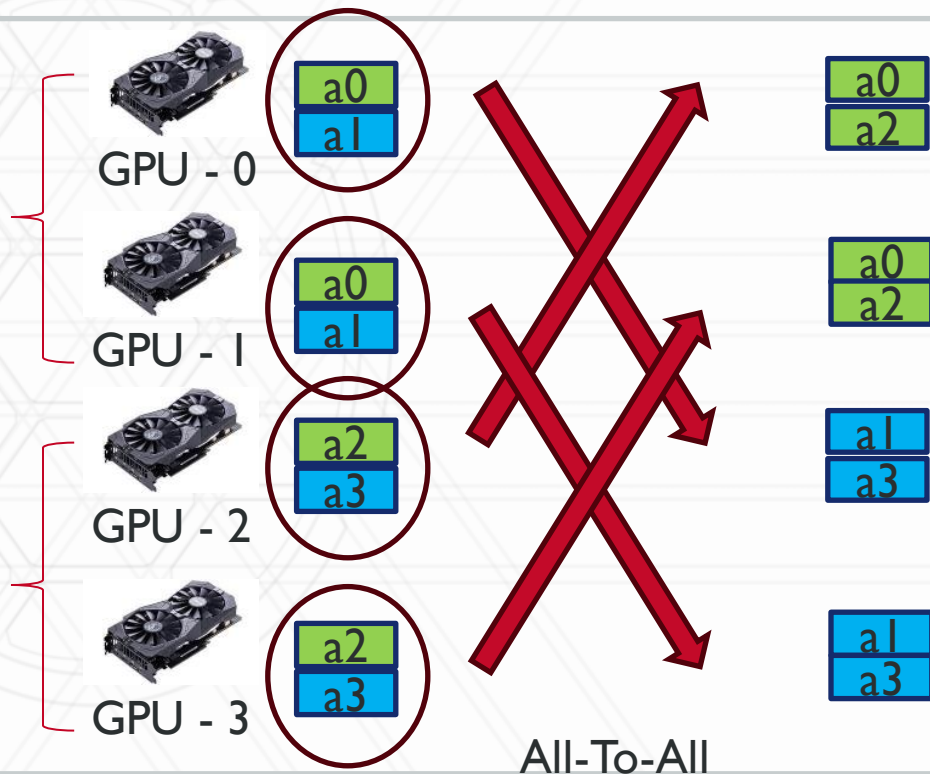
Largest Trainable MoE Models on Summit



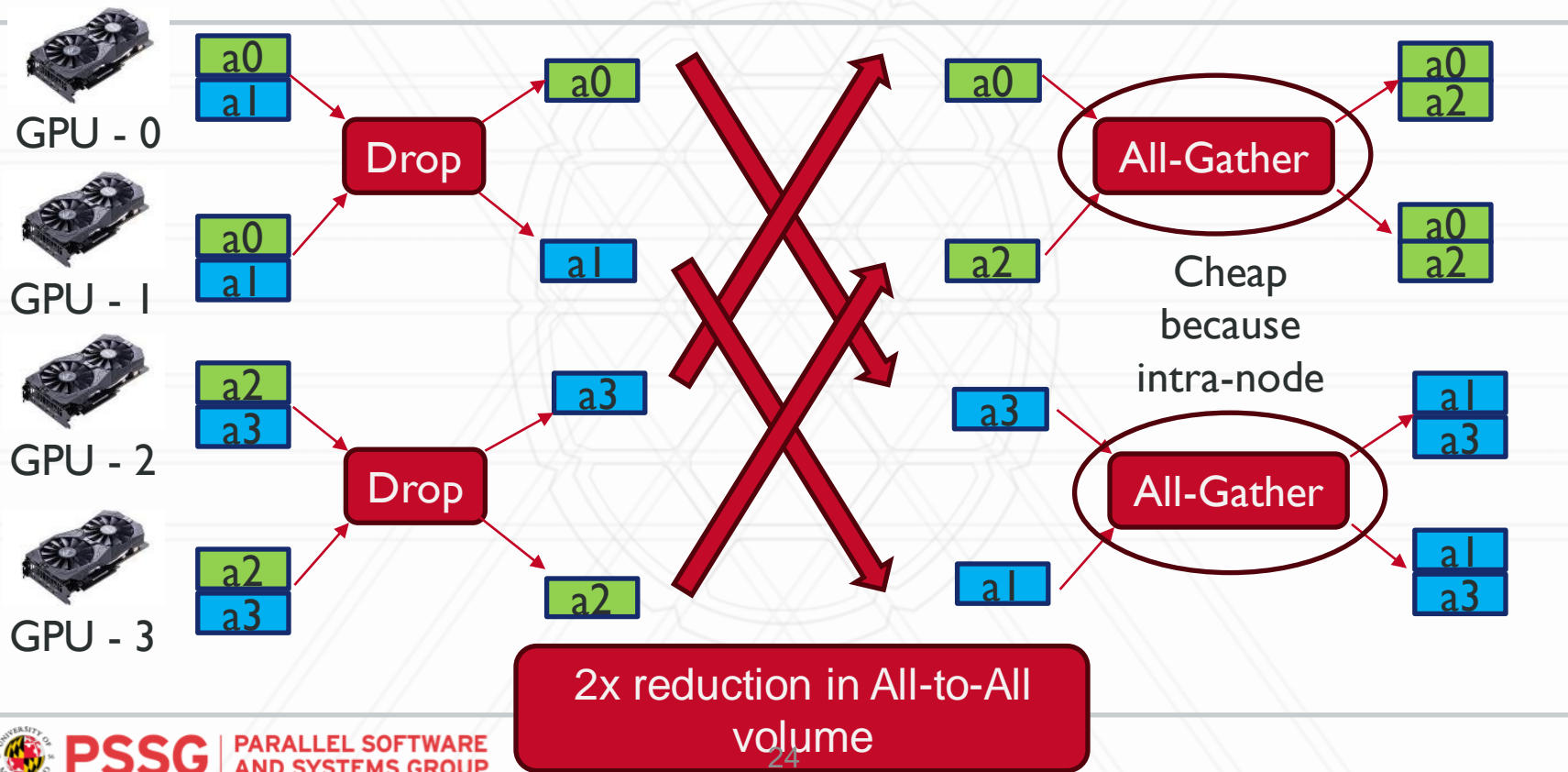
- Limit Number of experts to 128
- Limit tensor parallelism to a node.

Opt #1 Duplicate Token Dropping (DTD)

- Consider the first all-to-all.
- Tensor parallel GPUs communicate duplicate tokens.
- Remove this duplication to decrease All-To-All message sizes.



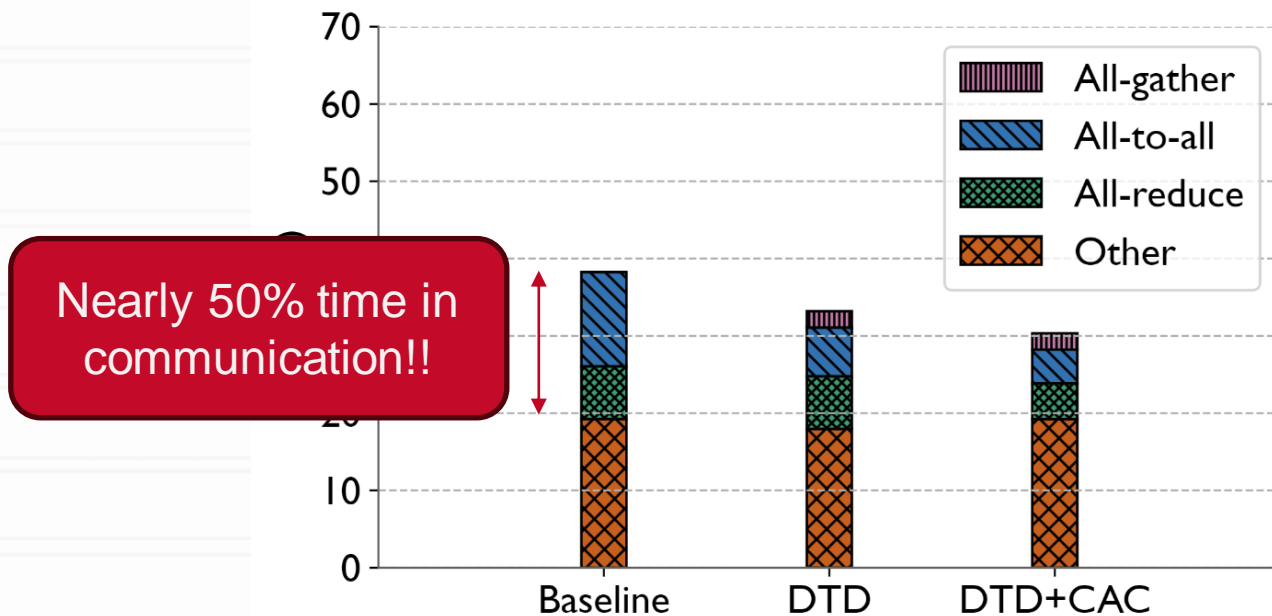
Opt #1 Duplicate Token Dropping



Opt #2 Communication-Aware Checkpointing (CAC)

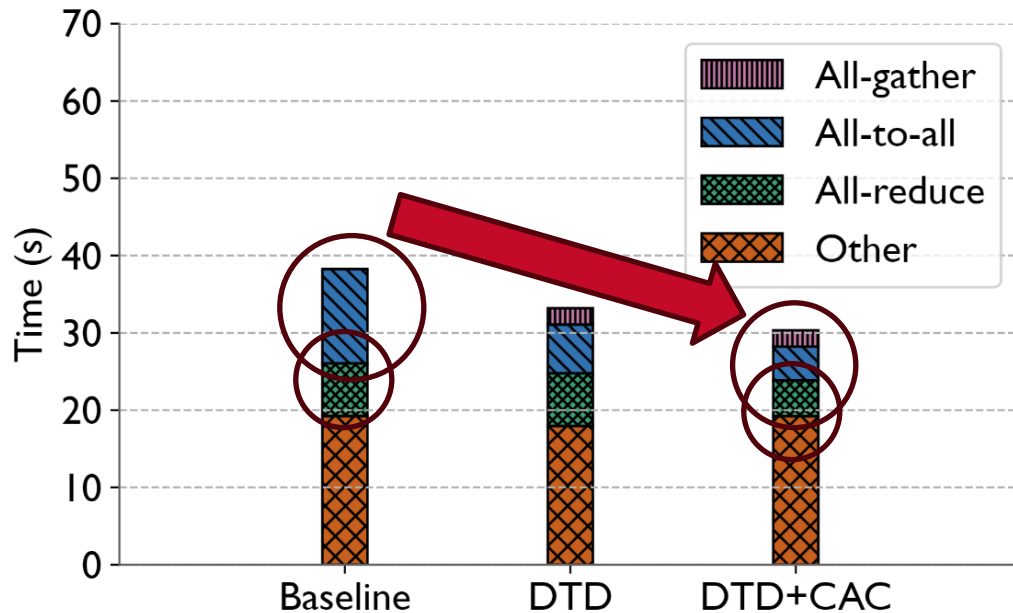
- Reduces number of all-to-all and all-reduce calls by 33 percent by utilizing marginally extra memory.
- More details in paper.

Results



Batch Time Profile of a 6.7B base model + 16 experts on 128 GPUs of Summit

Results

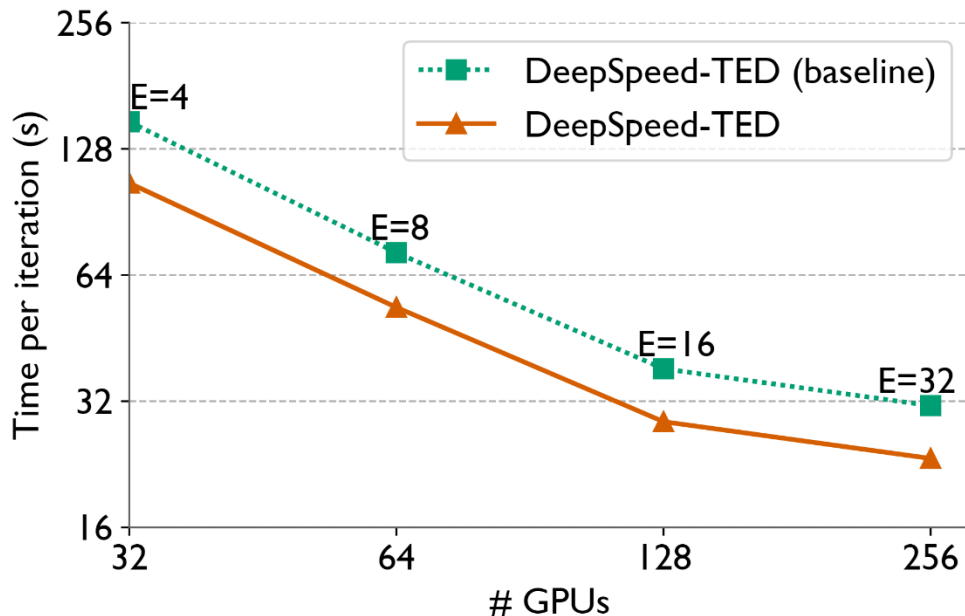


Overall 21%
Speedup

Batch Time Profile of a 6.7B base model + 16 experts on 128 GPUs of Summit

Results (Strong Scaling)

Strong Scaling of a 6.7B Base Model with Varying # Experts

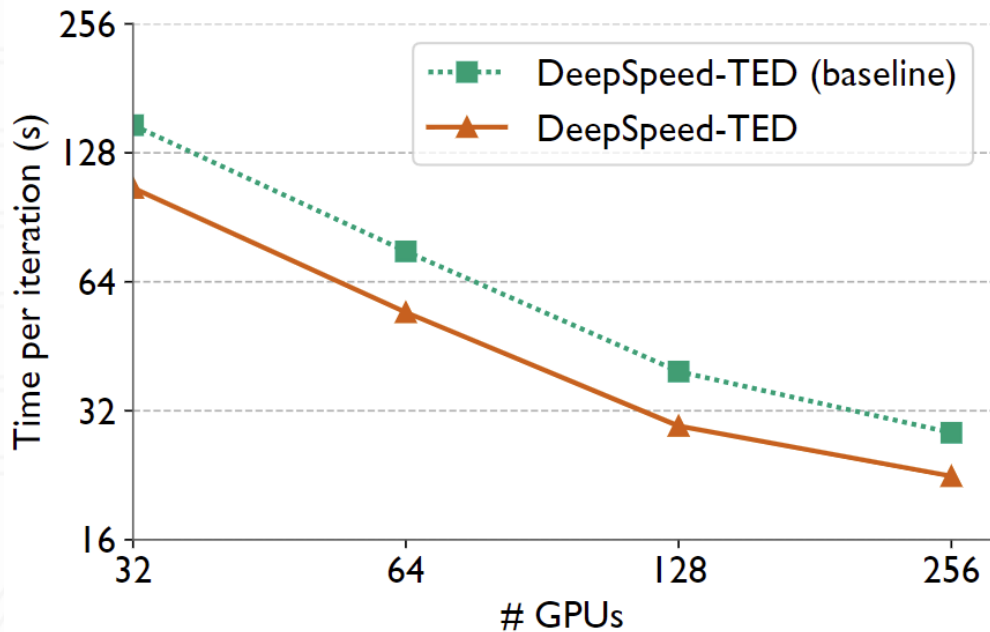


Machine - Summit

22-29%
speedups!

Results (Strong Scaling)

Strong Scaling for a 6.7B Base Model with 4 Experts



Machine - Summit



PSSG

PARALLEL SOFTWARE
AND SYSTEMS GROUP

Conclusion and Future Work

- Developed DeepSpeed-TED, a highly scalable parallel framework for training high quality MoEs with large base models.
- Presented a three-dimensional hybrid parallel method that supports MoEs with 4-8x larger models than the SoTA.
- Introduced communication optimizations that can achieve significant reductions in the collective communication times.
- As future work, we want to explore pipeline parallelism as a fourth dimension to scale to even larger base models.

Code

- Our work is integrated in DeepSpeed, a widely used open-source framework for parallel deep learning.
 - URL - <https://github.com/microsoft/DeepSpeed>

Bibliography

- [1] Using DeepSpeed and Megatron-LM to Train Megatron-LM Turing NLG 530B, A Large-Scale Generative Language Model, Smith et al., <https://arxiv.org/abs/2201.11990>
- [2] Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity, Fedus et al., <https://arxiv.org/abs/2101.03961>
- [3] Megatron-LM: Training Multi-Billion Parameter Language Models Using Model Parallelism, Shoeybi et al., <https://arxiv.org/abs/1909.08053>
- [4] DeepSpeed-MoE: Advancing Mixture-of-Experts Inference and Training to Power Next-Generation AI Scale, Rajbhandari et al., <https://arxiv.org/abs/2201.05596>
- [5] ZeRO: Memory Optimizations Toward Training Trillion Parameter Models, Rajbhandari et al., <https://arxiv.org/abs/1910.02054>



UNIVERSITY OF
MARYLAND

Siddharth Singh
ssingh37@umd.edu