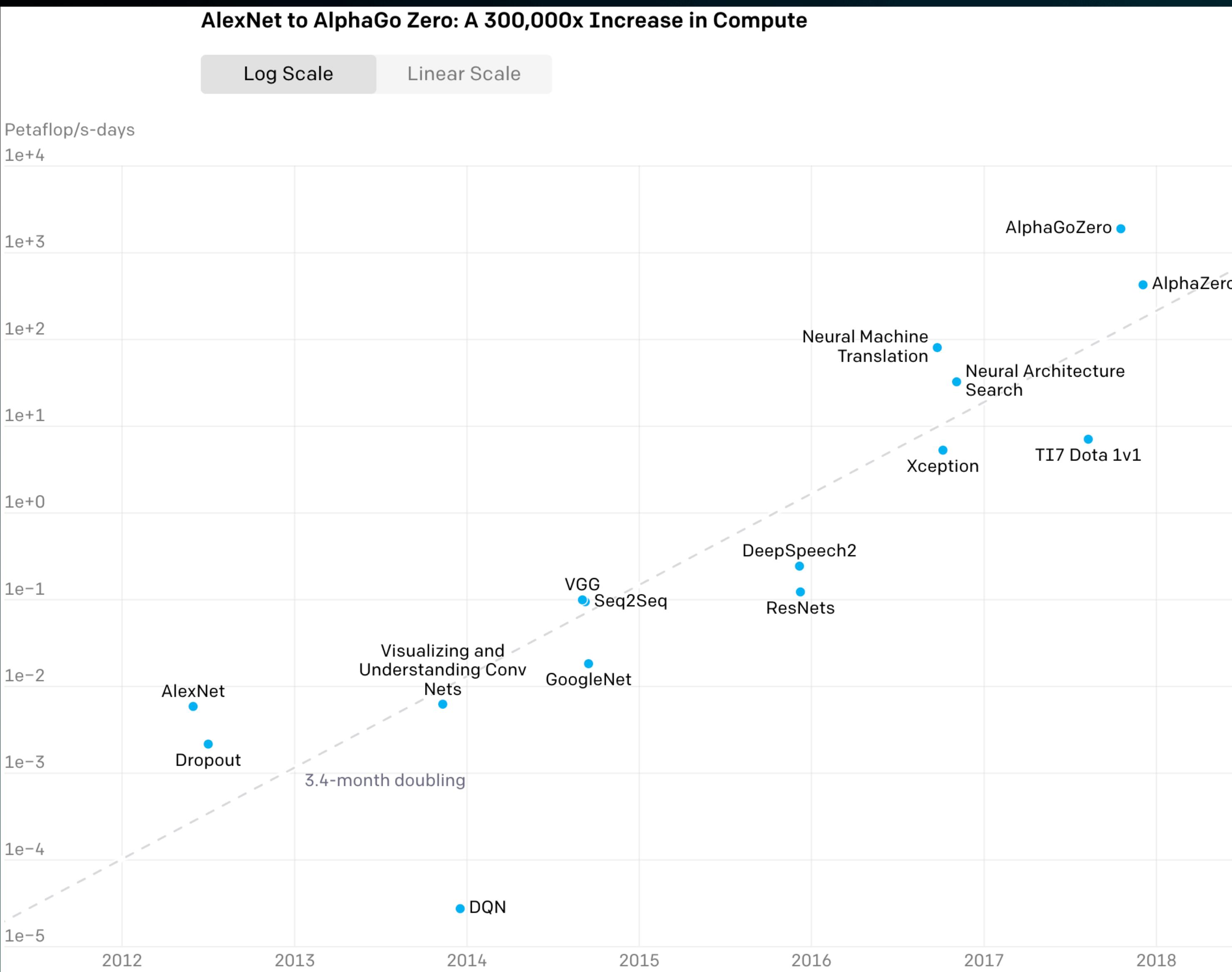


Optimizing Large-Scale GPT Models Pretraining

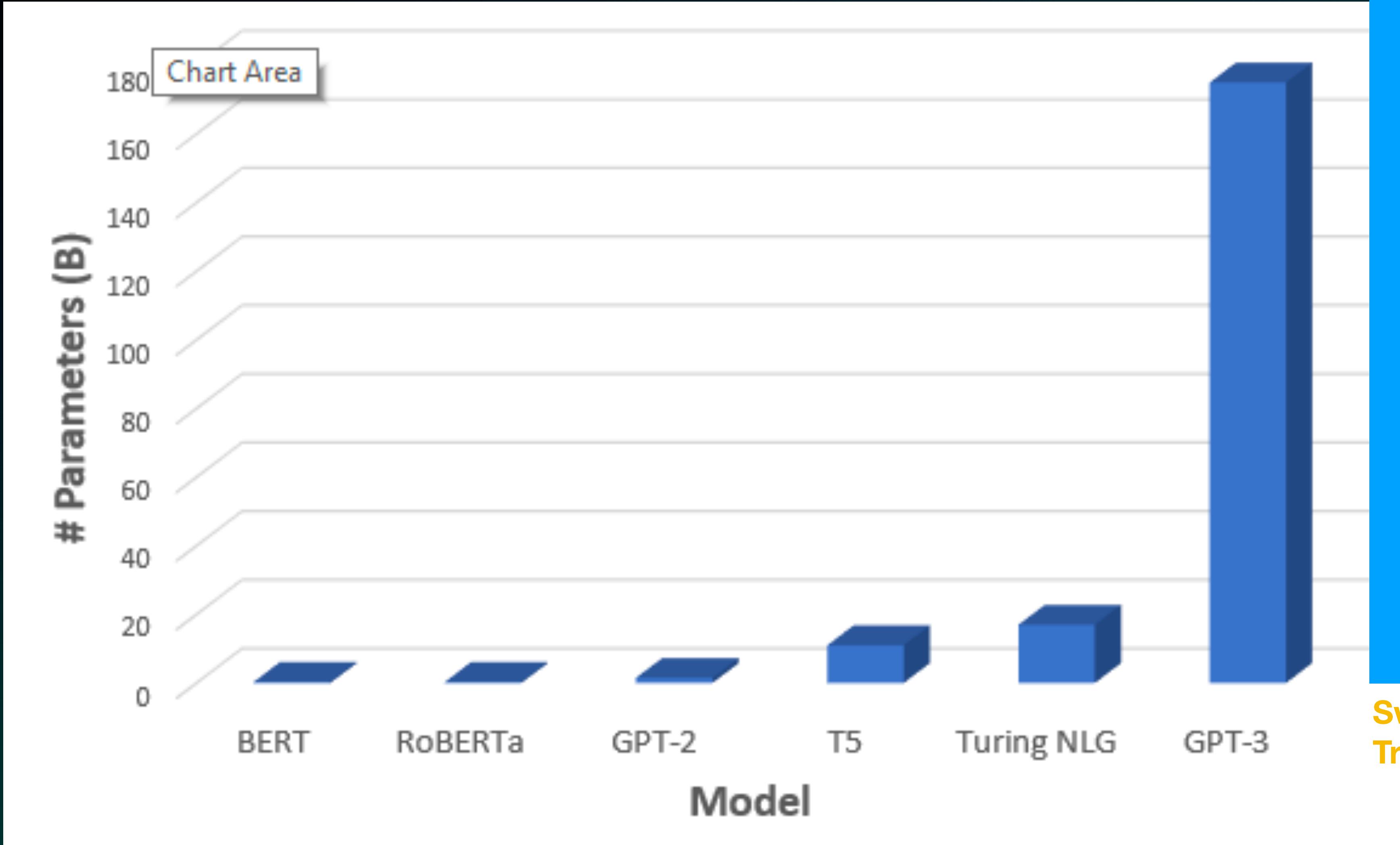
Oleh Shliazhko, NLP R&D, SberDevices

Large Scale Deep Learning NLP Models



1. Amount of SOTA compute doubling each 3.4 months
2. And this is ain't gonna stop anytime soon

Large Scale Deep Learning NLP Models



GPT-2 – 1.5 Billion weights
Megatron LM – 8.3 Billion
T5 – 11 Billion
Turing NLG – 17 Billion
GPT-3 – 175 Billion

Switch Transformer - 1.6 TRILLION 😱 😱 😱

Switch
Transformer

<https://arxiv.org/pdf/2101.03961.pdf>

Large Scale Deep Learning NLP Models

Large Scale Modeling limited by:

- Memory constraints
 - Model size
 - Batch size
 - Optimizer choice
- Time constraints
 - Lot of samples
 - Long step time

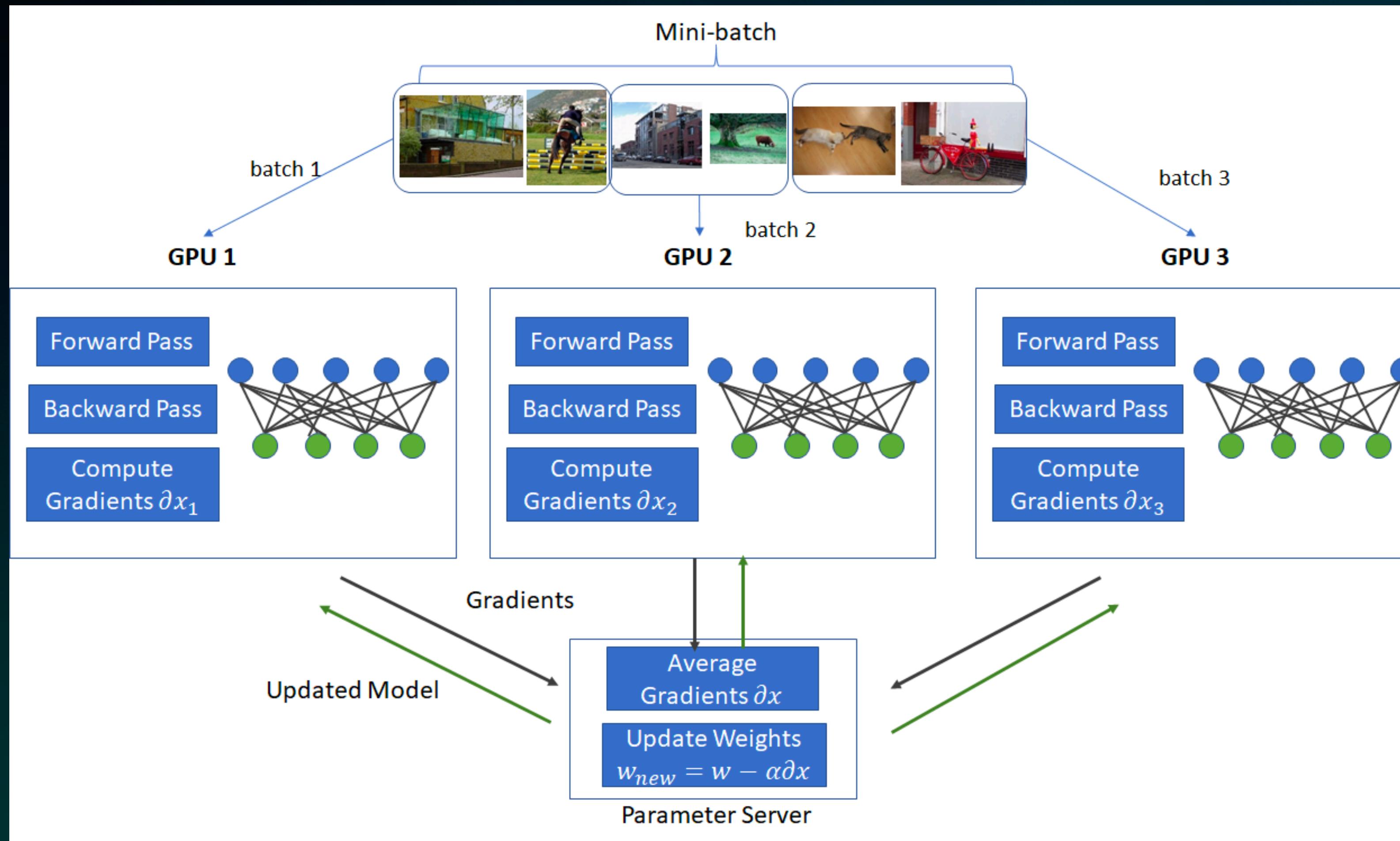
Our path to 13B GPT-3

1. GPT-2 small, 117M
2. GPT-2 small, + data parallel
3. GPT-2 medium, 345M
4. GPT-3-scale dataset
5. GPT-3 medium 345M
6. GPT-3 large, 750M, + deepspeed optimizer
7. GPT-3 sparse layer selection
8. GPT-3 XL, 1.3B, + checkpoint activations
9. GPT-3 13B, + model parallel
10.We need to go Deeper



Large Scale Deep Learning NLP Models

Usual workflow, distributed Data-Parallel Training



Memory constraints

Modern GPU Memory constraints:

Nvidia V100 - 32GB

Nvidia 2080TI - 11GB

.

.

Nvidia A100 - 80GB 

Weights → Bytes approximation

Training with FP32 Adam:

$$32 \text{ bit weights} + 32 \times 2 \text{ optimizer state} = \\ \underline{12 \text{ bytes}}$$

Training with FP16 Mixed precision:

$$16 \text{ weights} + 16 \times 2 + 32 \times 2 \text{ optimizer} = \\ \underline{14 \text{ bytes!}}$$

Training with pure FP16 training:

$$16 \text{ weights} + 16 \times 2 \text{ optimizer} = 6 \text{ bytes}$$

*But not available out of the box 😞

Upper bound model size

- + some memory for input batches (megabytes)
- + some memory for activations (few gigabytes)

Practical upper bound of GPT model size per GPU:

- V100 - 2 billion weights
- 2080TI - 700 million weights

Data parallelism is not enough.

Ways to go:

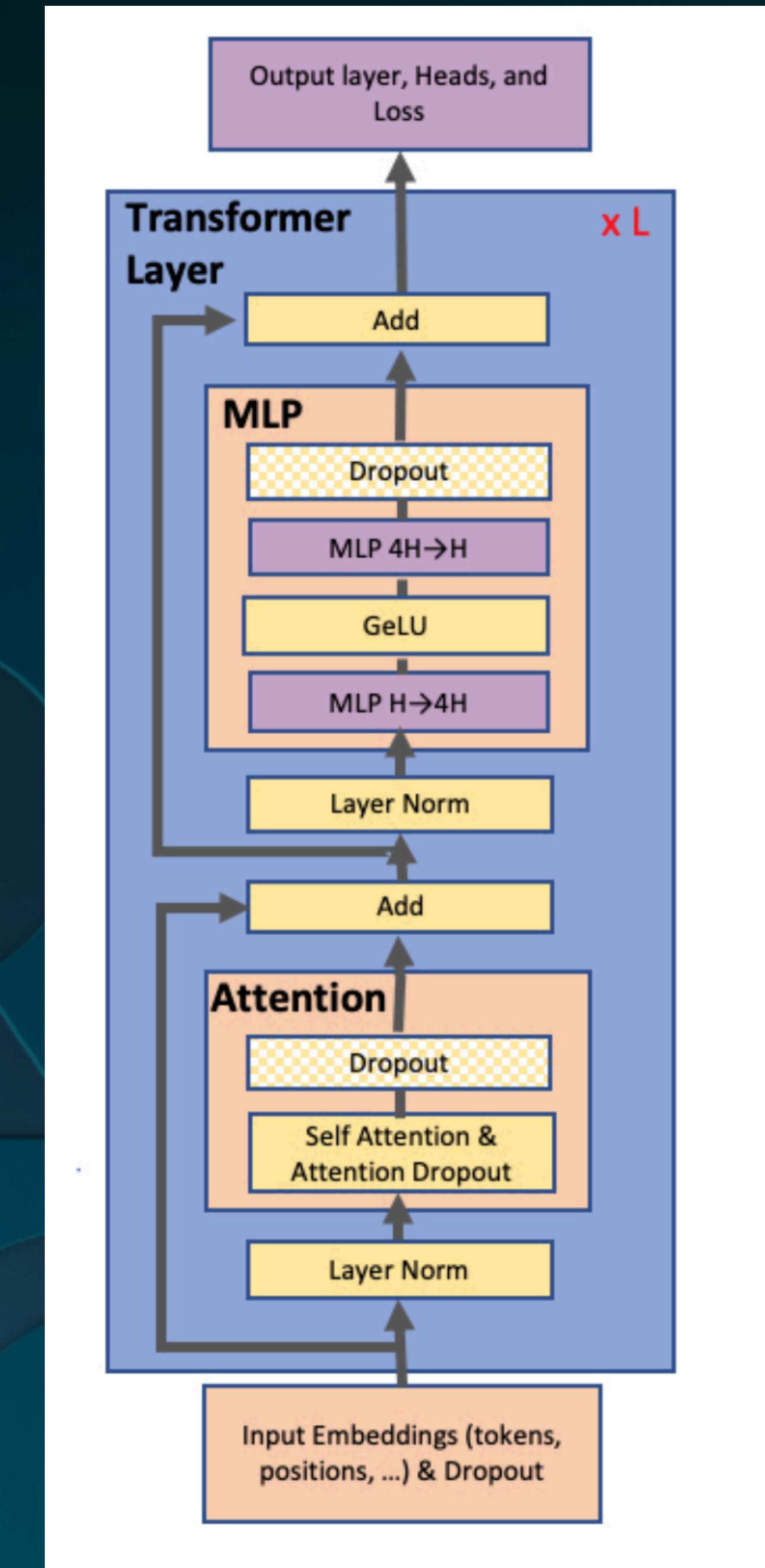
- Split the model between GPUs
- Split the optimizer between GPUs
- Recompute activations instead of store

Split the model between GPUs

Model parallelism

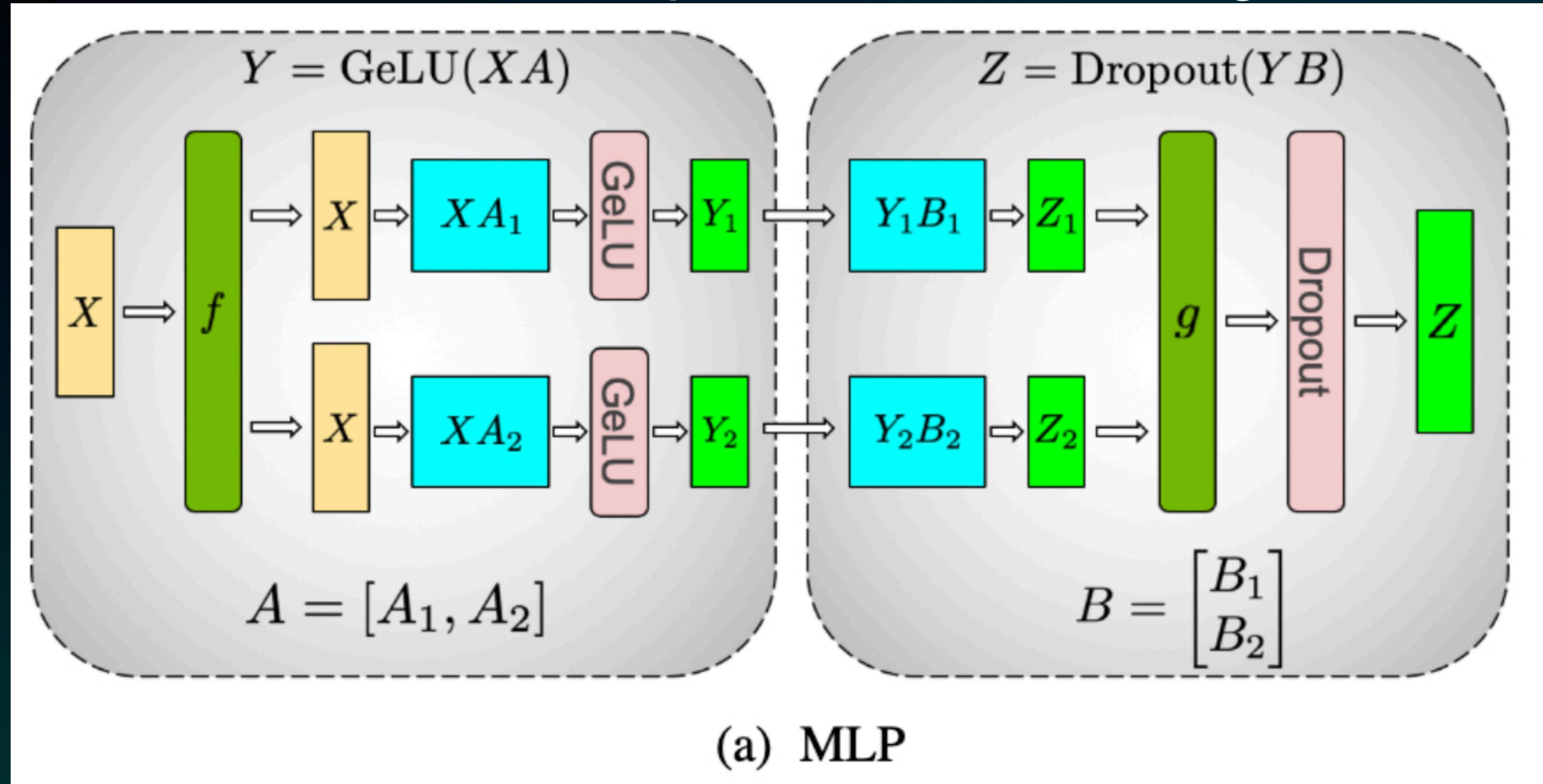
- Pipeline (layer-wise) parallelism
- Distributed tensor computation

Transformer Reminder



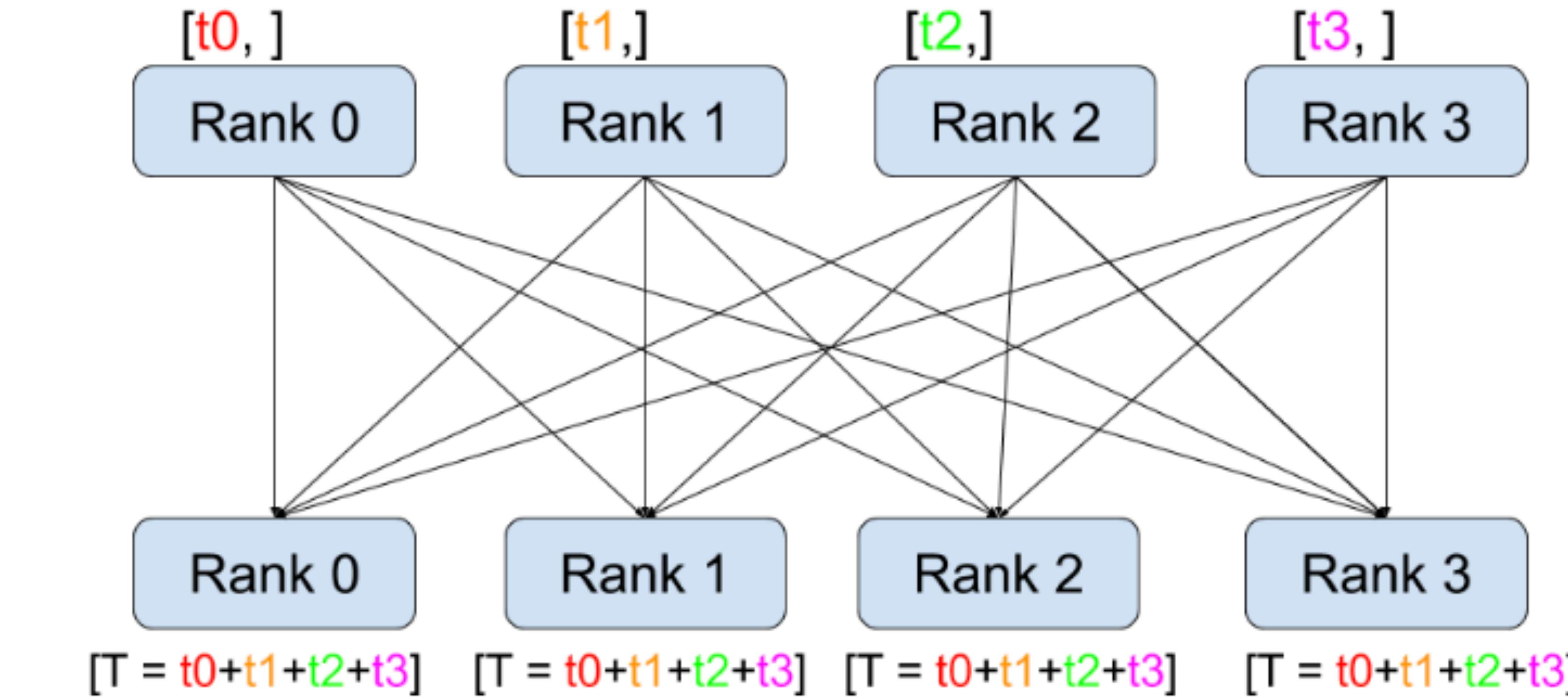
Split the model between GPUs

Distributed tensor computation. Nvidia Megatron



```
class f(torch.autograd.Function):
    def forward(ctx, x):
        return x
    def backward(ctx, gradient):
        all_reduce(gradient)
        return gradient
```

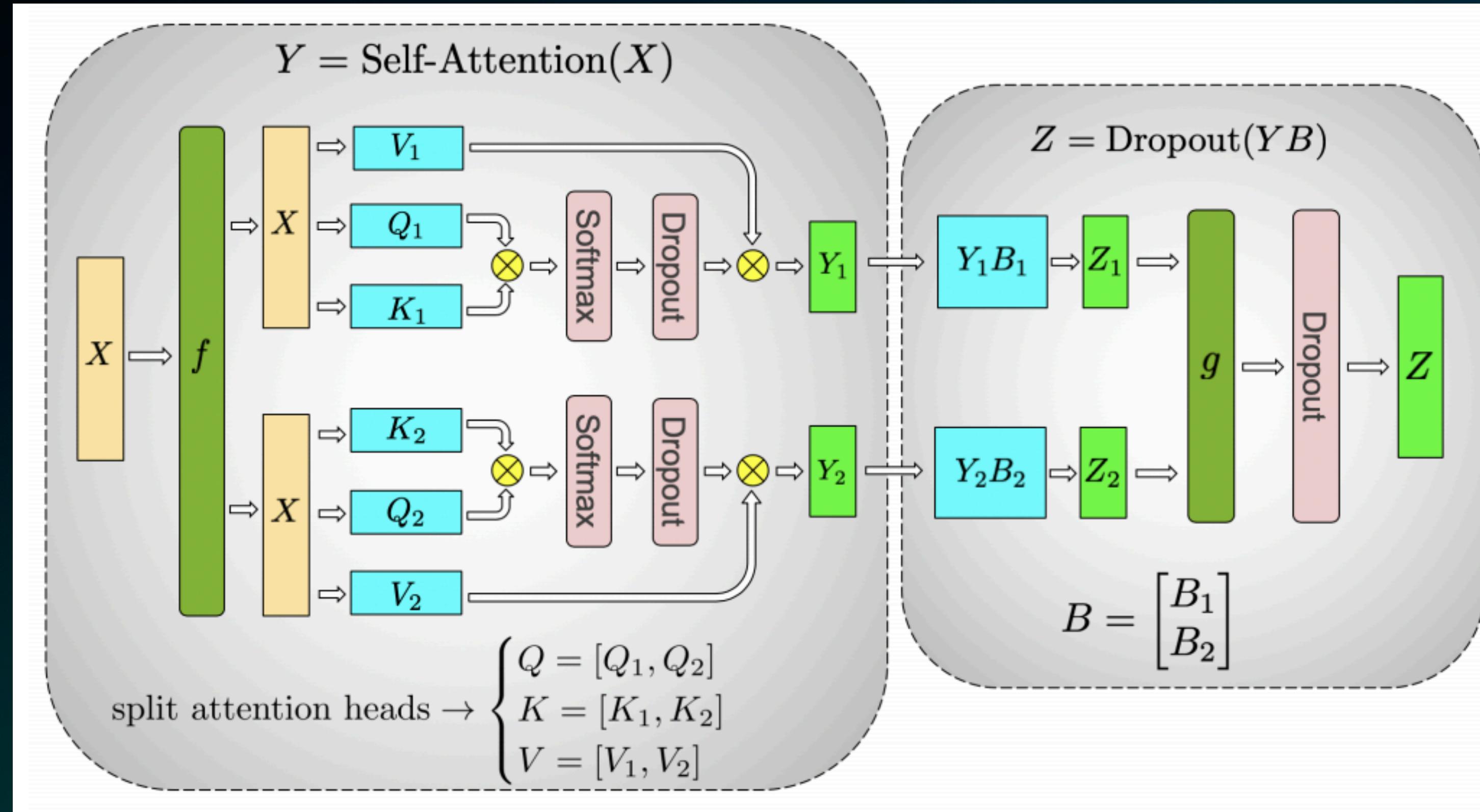
All-Reduce Reminder



All-Reduce

Split the model between GPUs

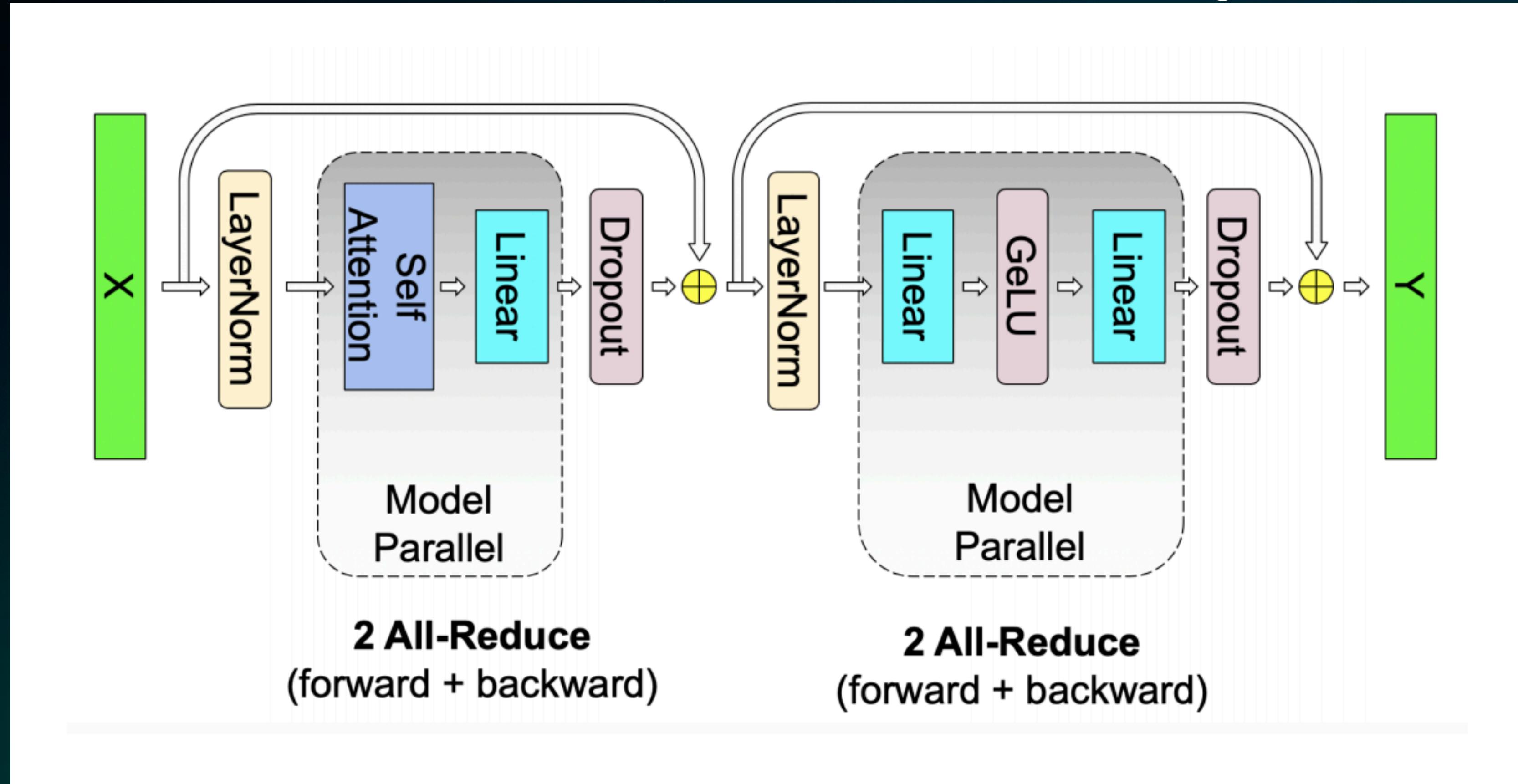
Distributed tensor computation. Nvidia Megatron



Partitioning the Q, K, V in a column parallel fashion such that the matrix multiply corresponding to each attention head is done locally on one GPU.

Split the model between GPUs

Distributed tensor computation. Nvidia Megatron



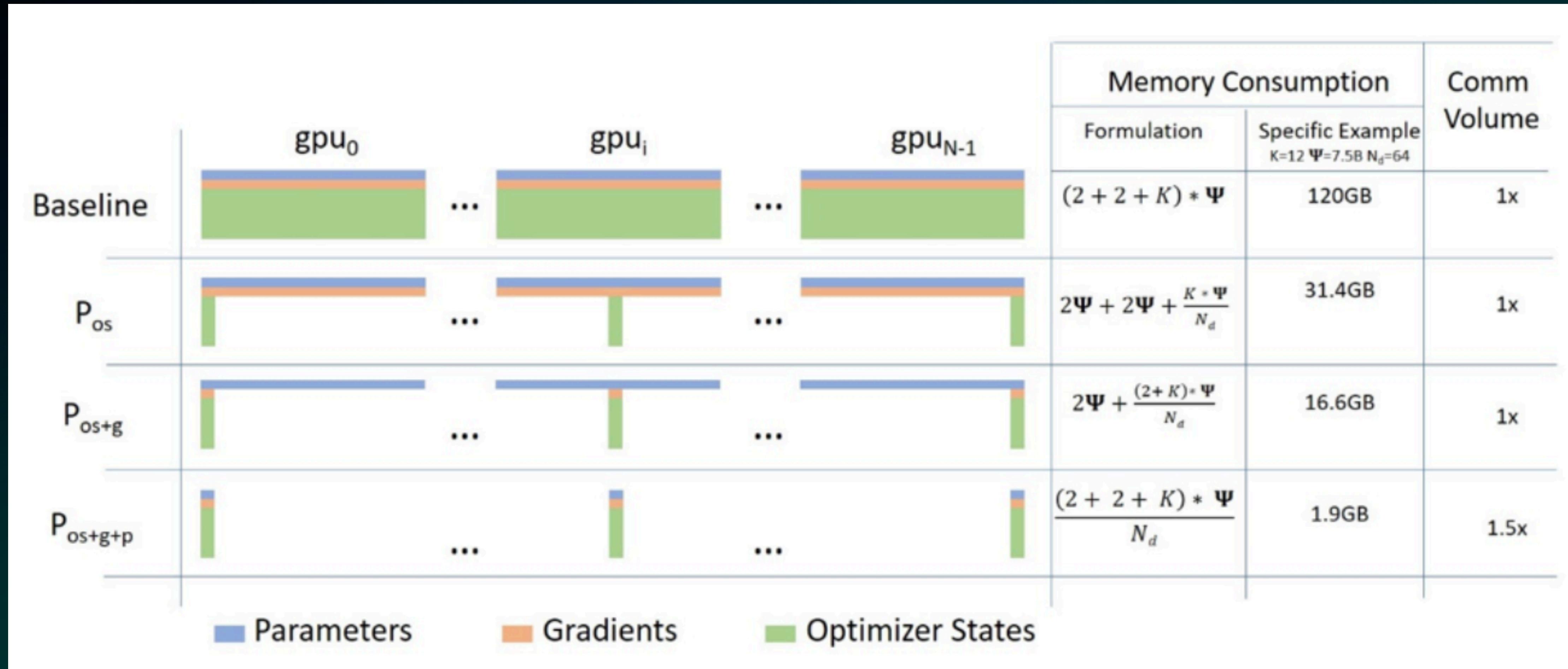
Split the optimizer between GPUs

Microsoft ZERO

1. Optimizer State Partitioning (Pos) – 4x memory reduction, same communication volume as data parallelism. Update locally, then all_gather.
2. Add Gradient Partitioning (Pos+g) – 8x memory reduction, same communication volume as data parallelism
3. Add Parameter Partitioning (Pos+g+p) – Memory reduction is linear with data parallelism degree Nd.

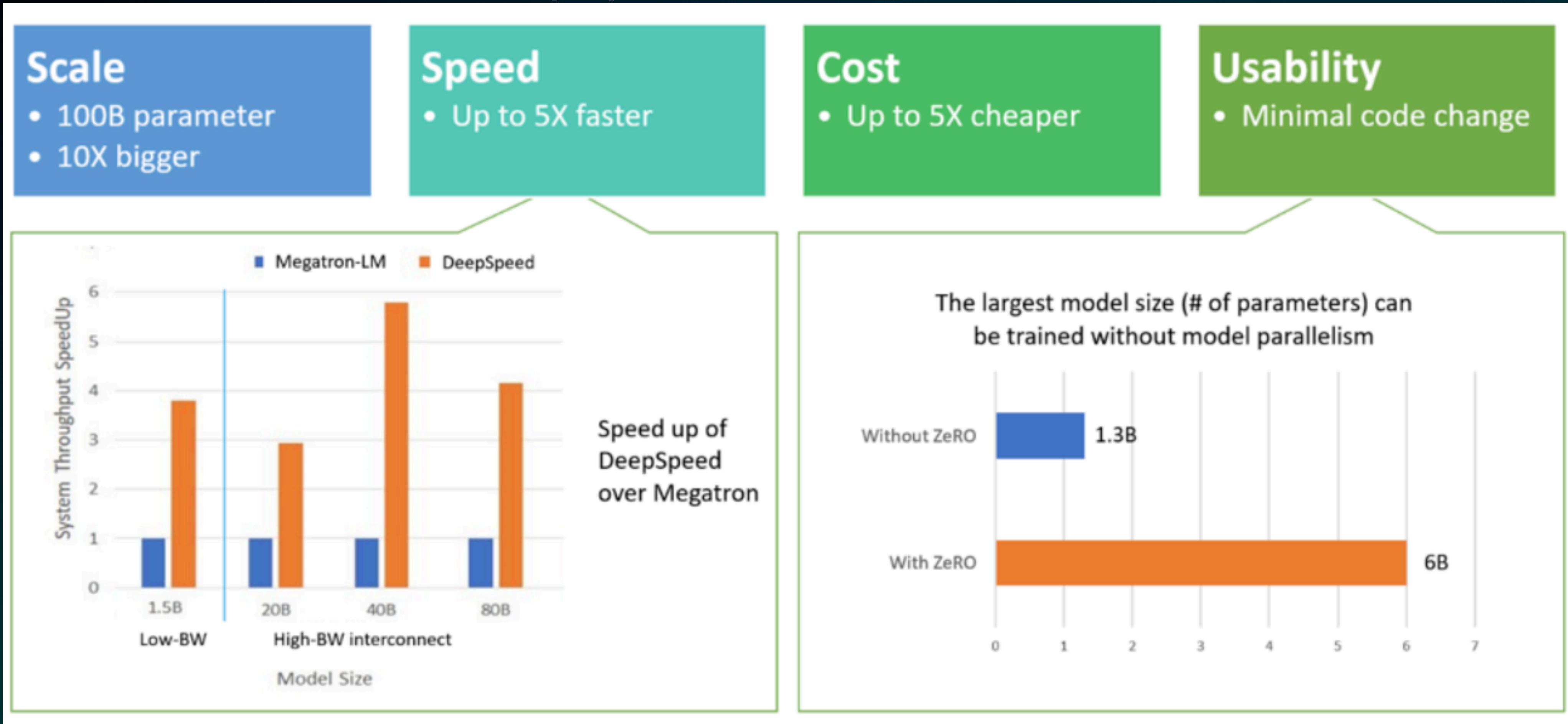
Split the optimizer between GPUs

Microsoft ZERO



Split the optimizer between GPUs

Microsoft ZERO & Deepspeed



Activation Checkpointing

- saves memory
- store only the activations of each N th layer
- recompute others during backprop

<https://arxiv.org/pdf/1904.10631.pdf>

Time constraints

Self-attention has quadratic ($N \times N$) complexity.

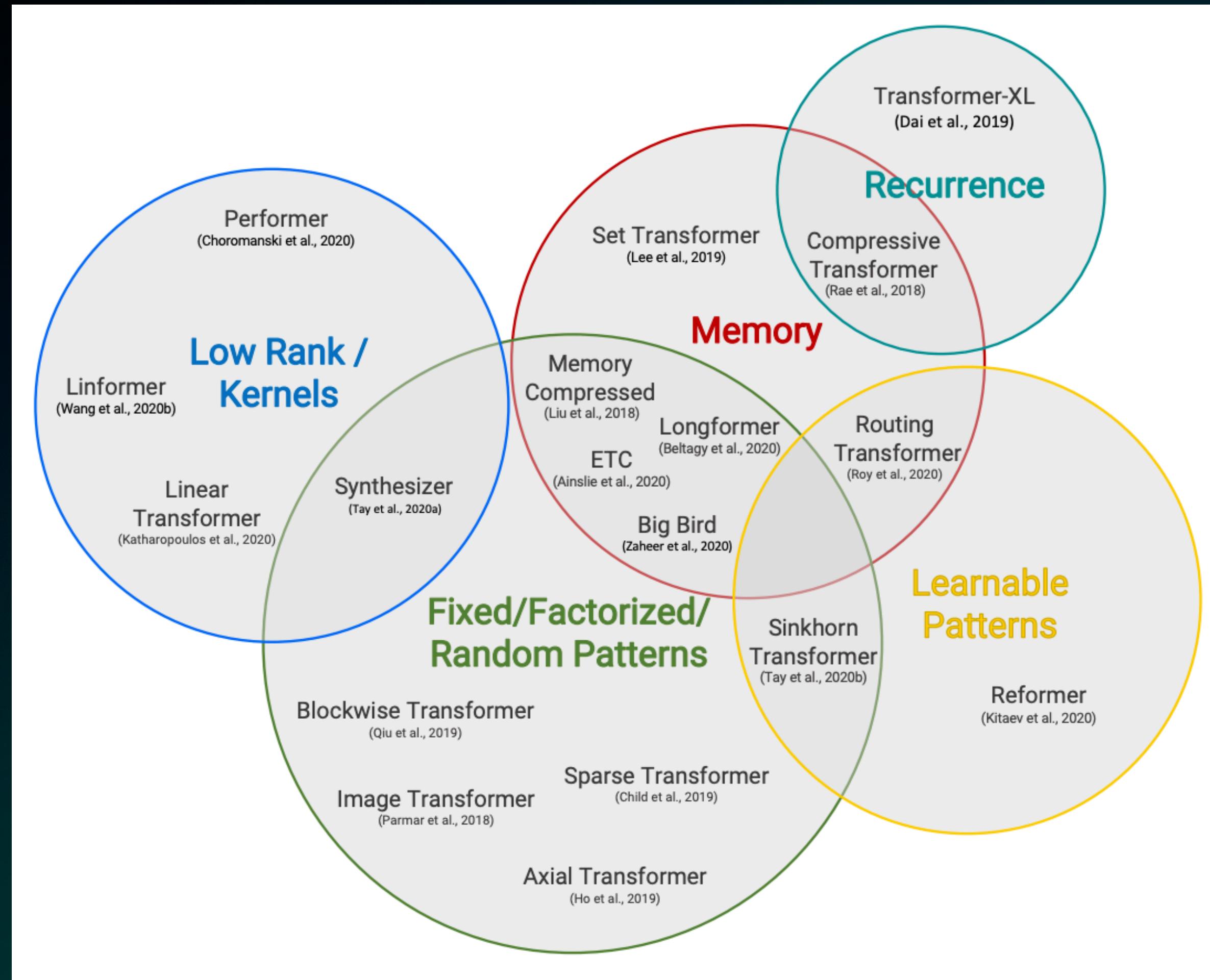
Most of the compute time of transformer model spent in attention matrix computation.

Attention as weights of components - only few biggest matter most.

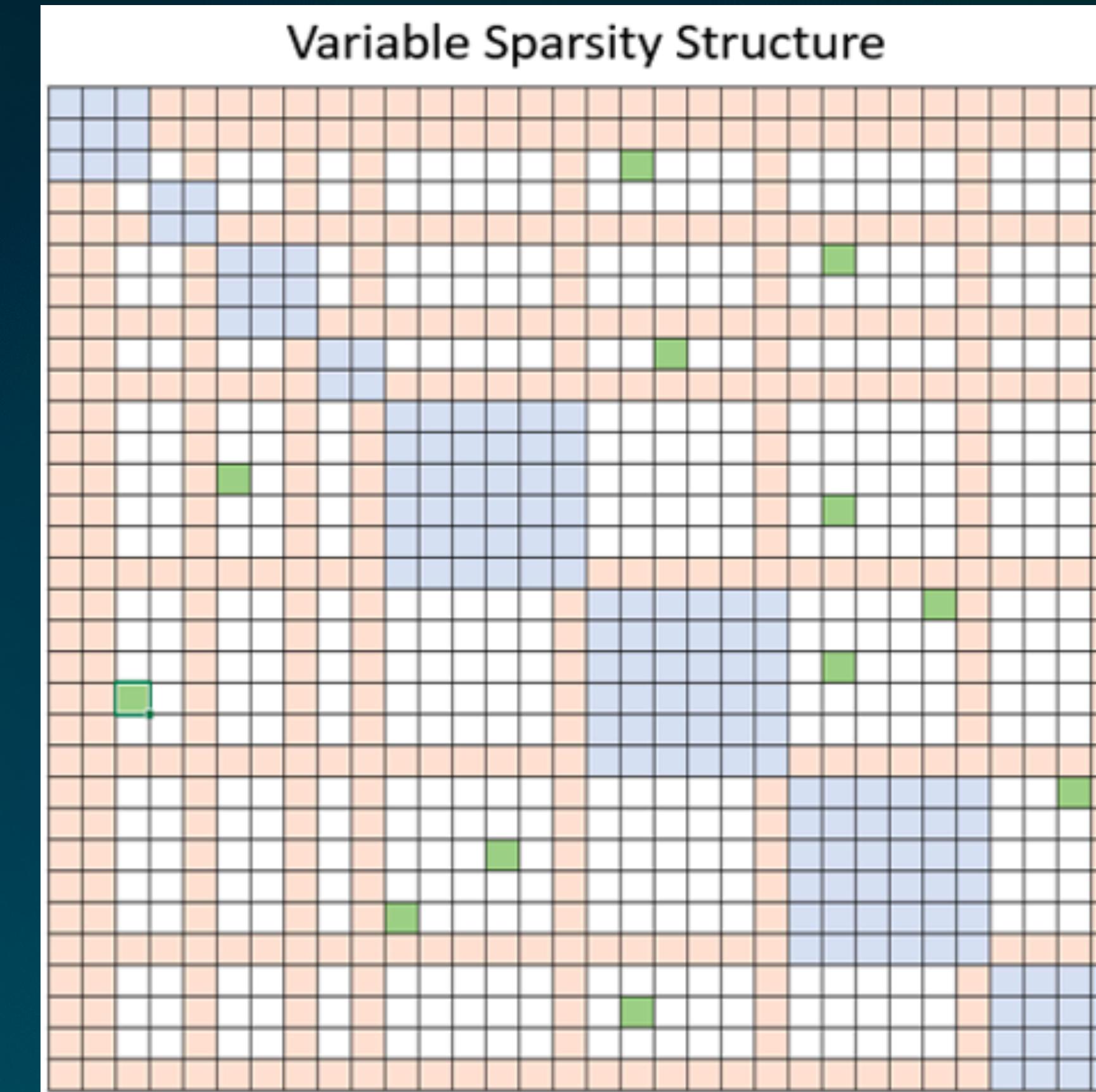
So, lets approximate this matrix with something faster.

<https://arxiv.org/pdf/2009.06732.pdf>

Time constraints



Exploit inherent attention matrix sparsity



<https://arxiv.org/pdf/2009.06732.pdf>

Thank you!

