

# Efficient LLMs

Yatin Nandwani  
Research Scientist, IBM Research



Semester 1,  
2025-2026

Large Language Models: Introduction and Recent Advances

ELL881 · AIL821

# IBM Research, India

## Conversational-AI



**Yatin  
Nandwani**



**Sonam  
Mishra**



**Dinesh  
Khandelwal**



**Gaurav  
Pandey**



**Vineet  
Kumar**



**Meghanadh  
Pulivarthy**



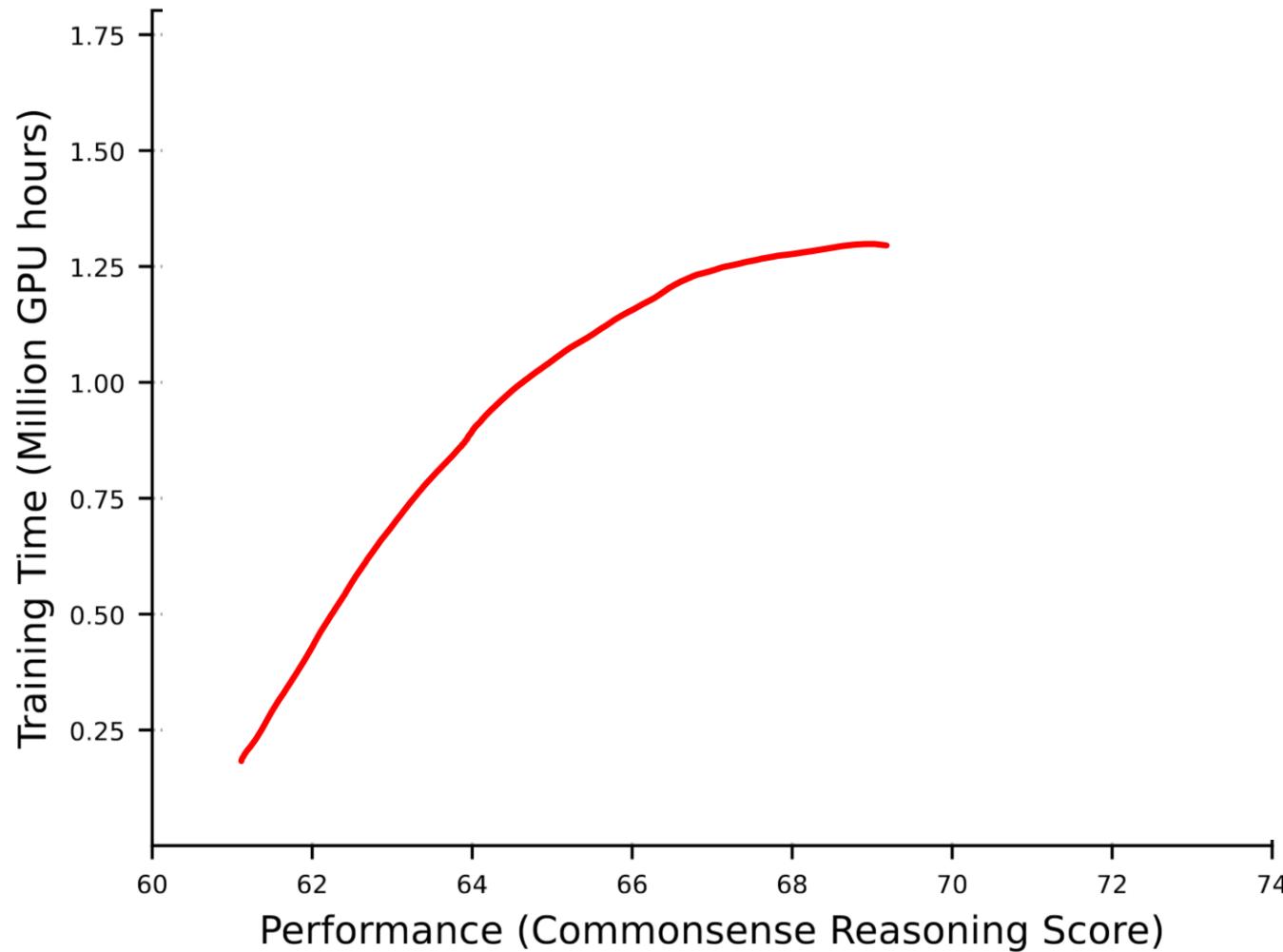
**Kushagra  
Bhushan**



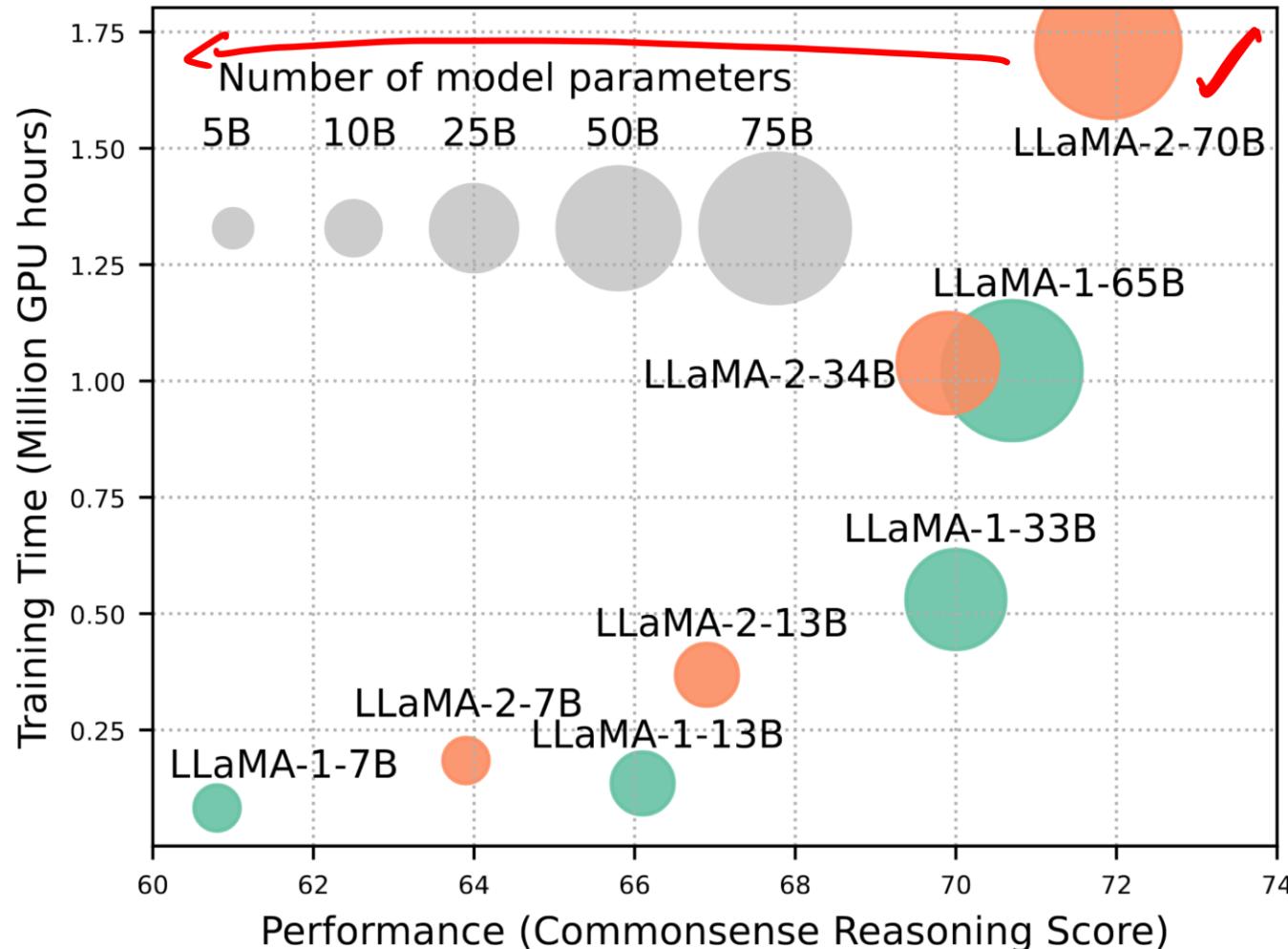
**Dinesh  
Raghu**

**Sachindra  
Joshi**

# Training Resources vs Performance



# Training Resources vs Performance



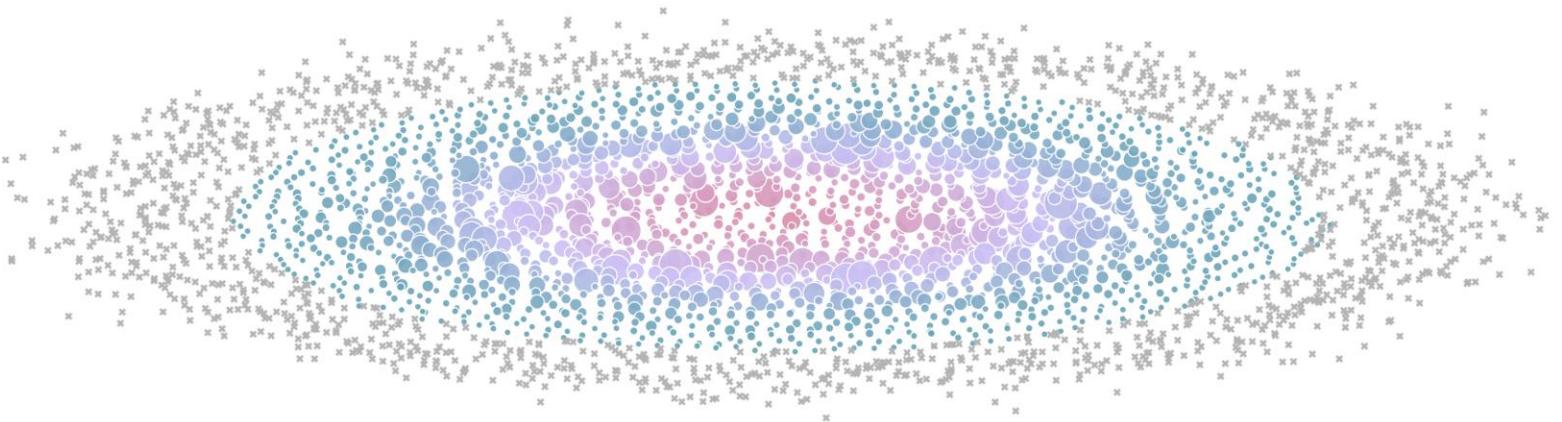
- Based on Nvidia A100 80GB GPU

<https://huggingface.co/spaces/optimum/llm-perf-leaderboard>

# Efficient LLMs

- How to scale training?
  - *Data Parallelism*
  - *Tensor Parallelism*
  - *Context Parallelism*
  - *Pipeline Parallelism*

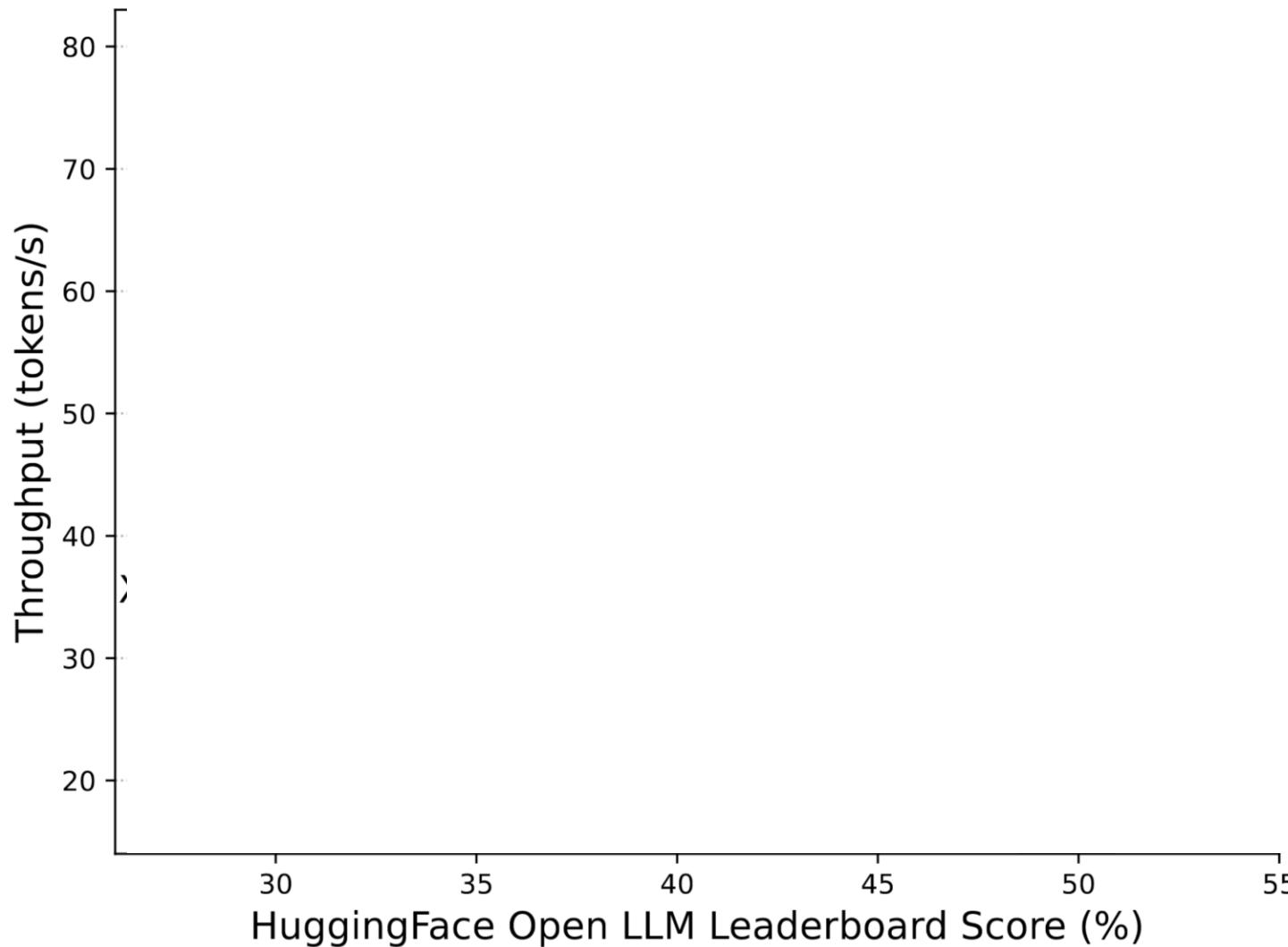
## The Ultra-Scale Playbook: Training LLMs on GPU Clusters



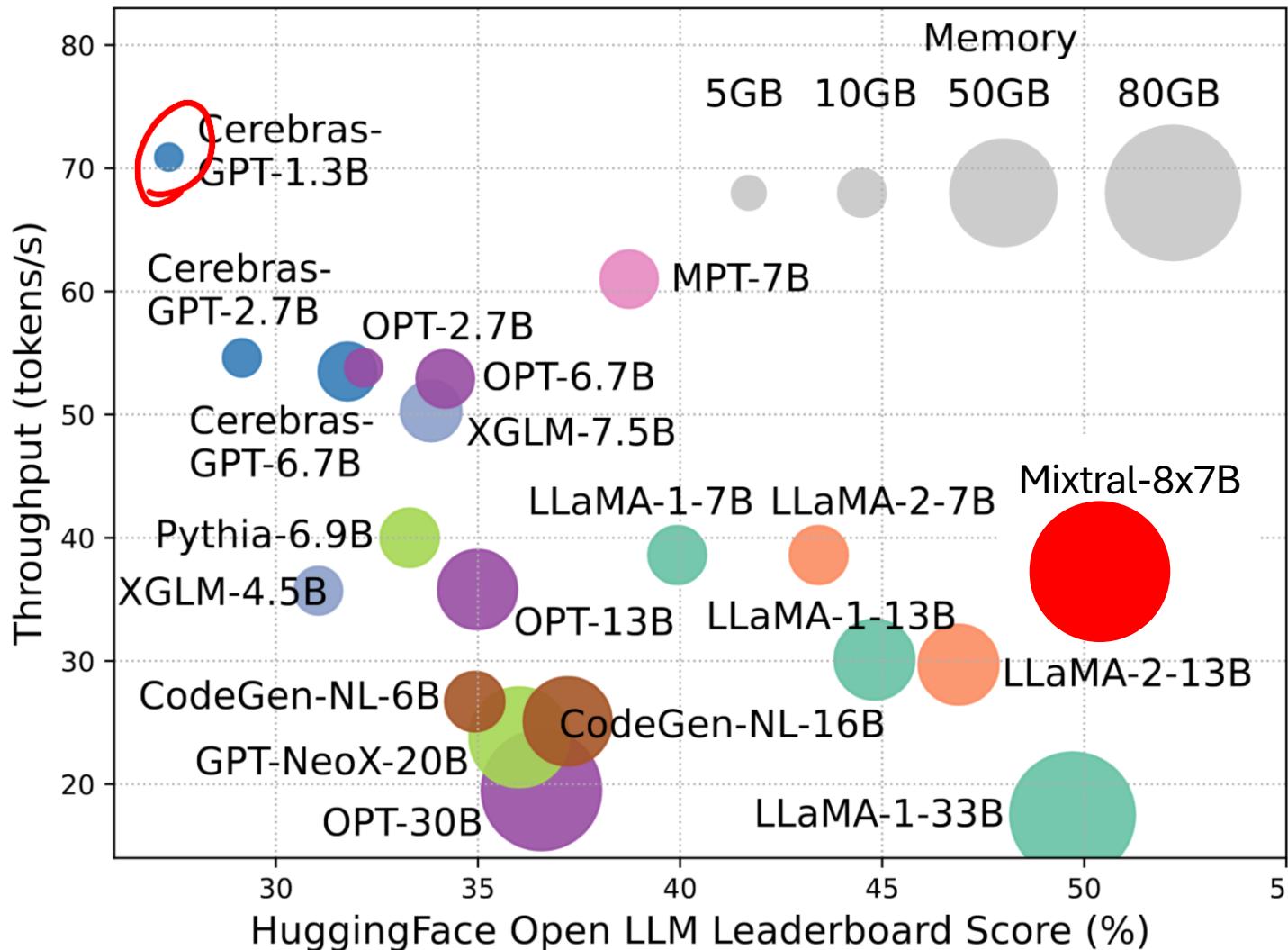
We ran over 4,000 scaling experiments on up to 512 GPUs and measured throughput (size of markers) and GPU utilization (color of markers). Note that both are normalized per model size in this visualization.



# Inference Throughput vs Performance

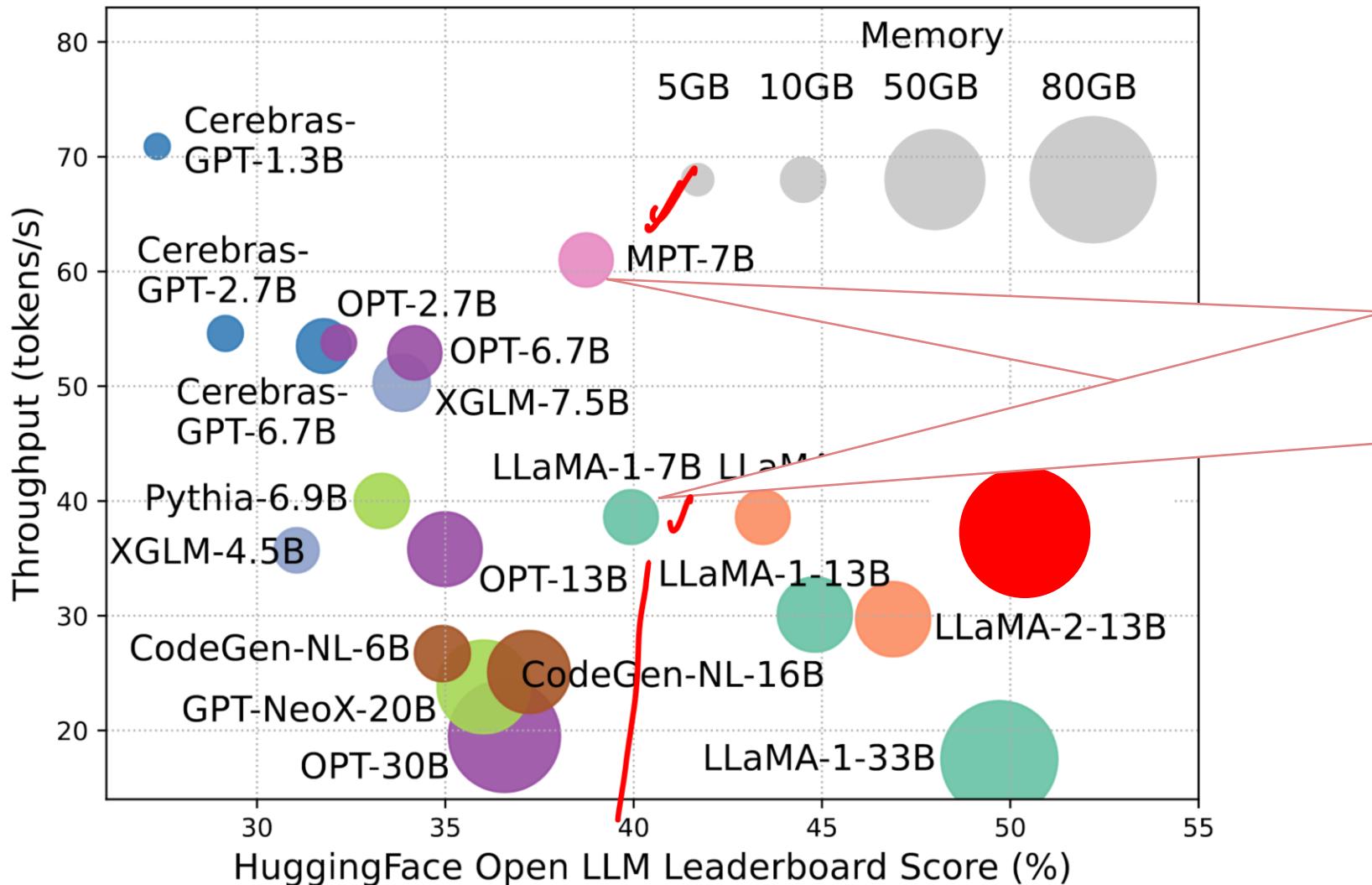


# Inference Throughput vs Performance



- On Nvidia A100 80GB GPU;
- 16-bit quantized
- Batch Size - 1
- Prompt size of 256
- Generating 1000 tokens

# Inference Throughput vs Performance



- Similar performance, different throughput! How?
- Efficient implementation –
  - Fused kernel for attention

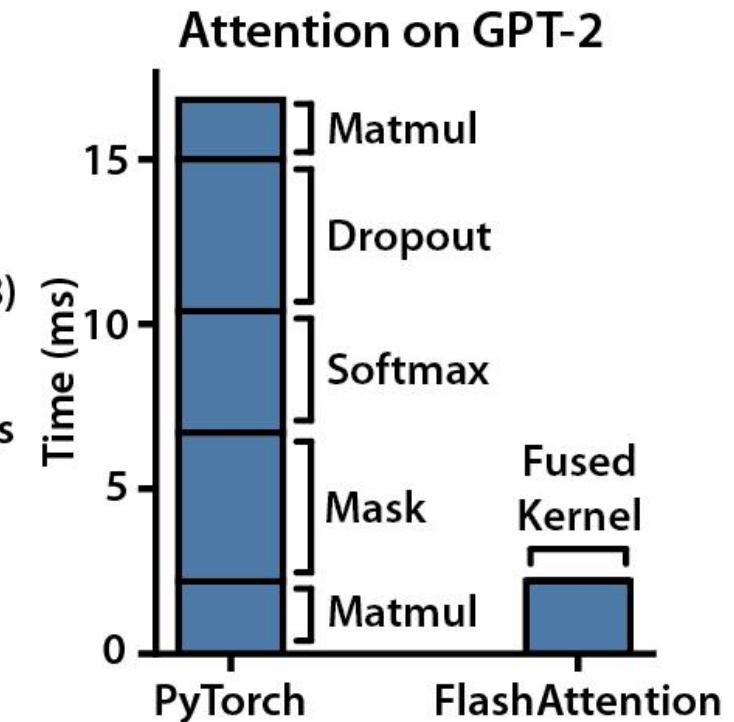
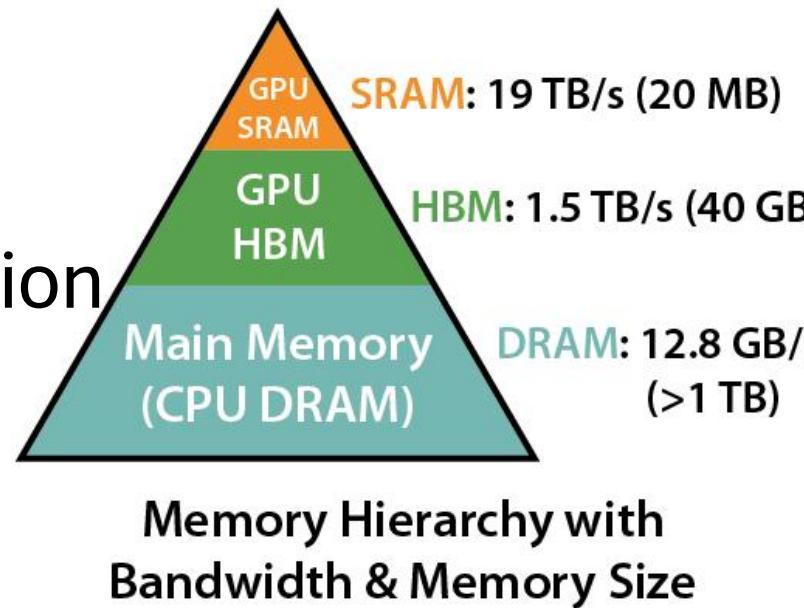
# Efficient LLMs

- How to scale training?

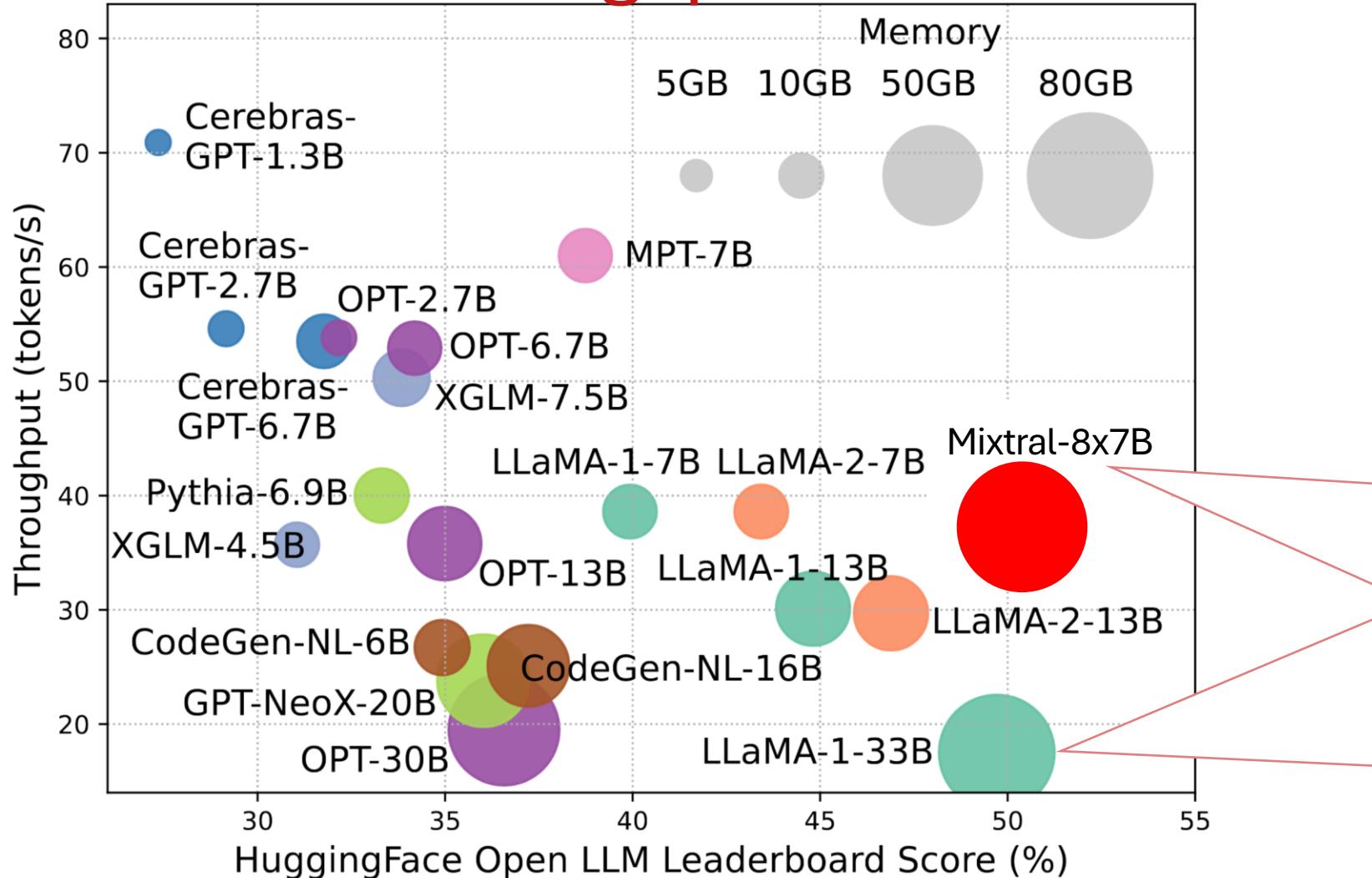
*Parallelism ...*

- Efficient implementation

- Flash Attention
- Paged Attention



# Inference Throughput vs Performance

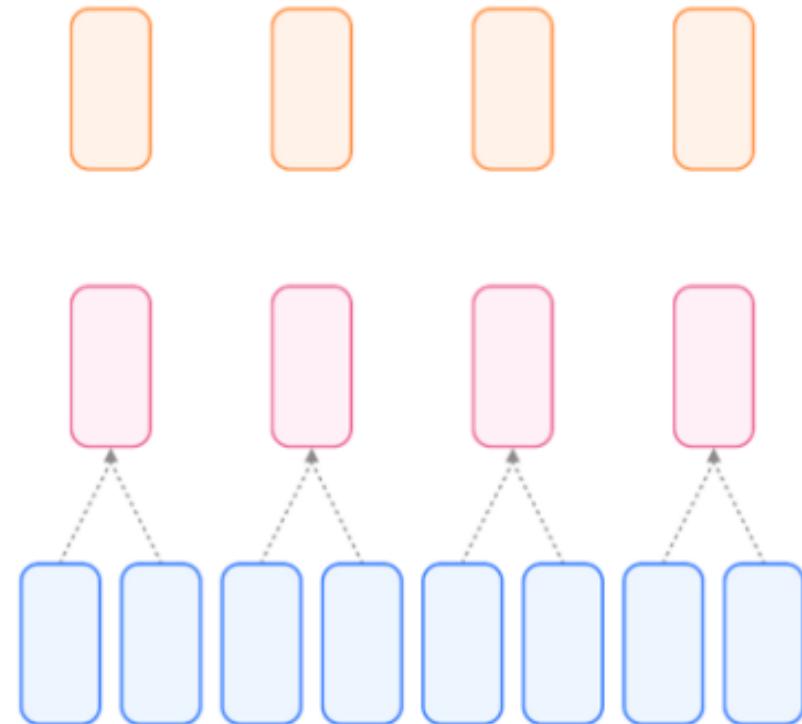


- Similar performance, different throughput! How?
- Efficient design –
  - Grouped Query Attention
  - MoE

# Efficient LLMs

- How to scale training?  
*Parallelism ...*
- Efficient implementation
  - *Fused kernels ...*
- Efficient design
  - Grouped Query Attention
  - Mixture of Experts

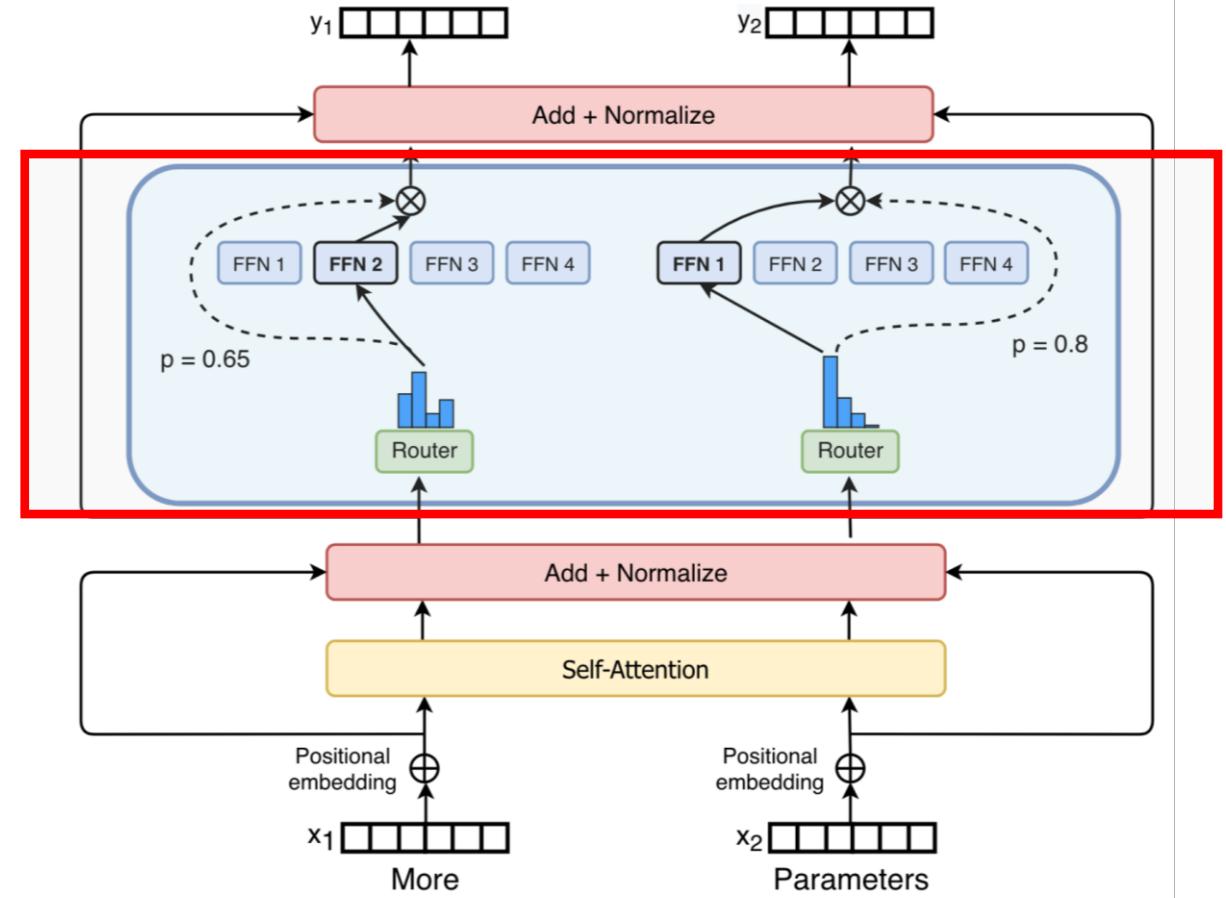
Grouped-query



# Efficient LLMs

- How to scale training?  
*Parallelism ...*
- Efficient implementation
  - *Fused kernels ...*
- Efficient design
  - Grouped Query Attention
  - Mixture of Experts

## Mixture of Experts



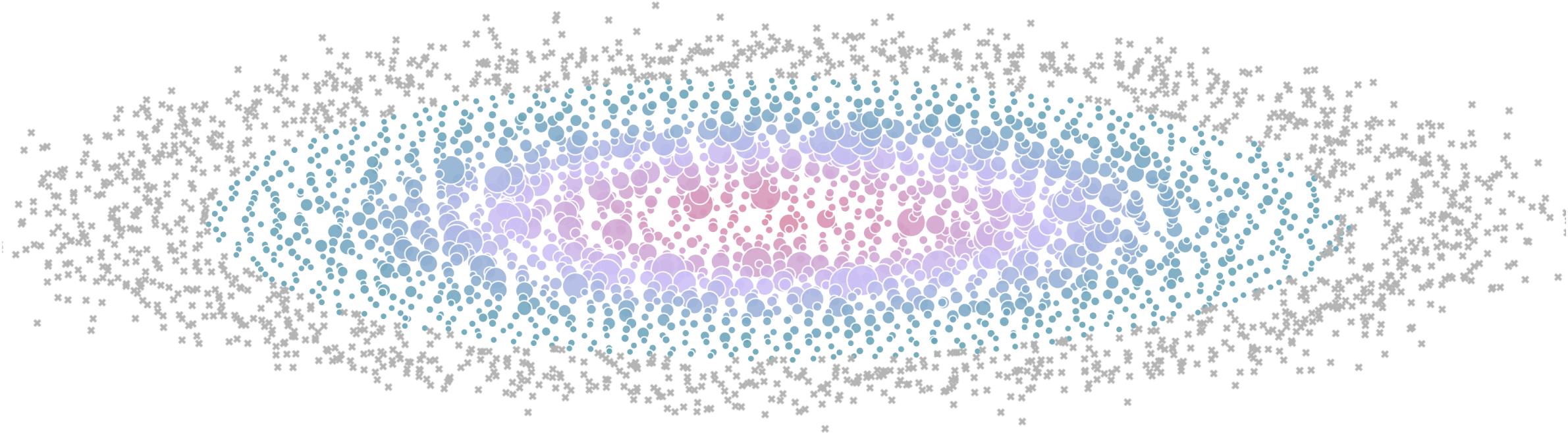
# Efficient LLMs

- How to scale training?  
*Parallelism ...*
- Efficient implementation
  - *Fused kernels ...*
- Efficient design
  - *Grouped Query Attention*
  - *Mixture of Experts*



# The Ultra-Scale Playbook: Training LLMs on GPU Clusters

<https://huggingface.co/spaces/nanotron/ultrascale-playbook>

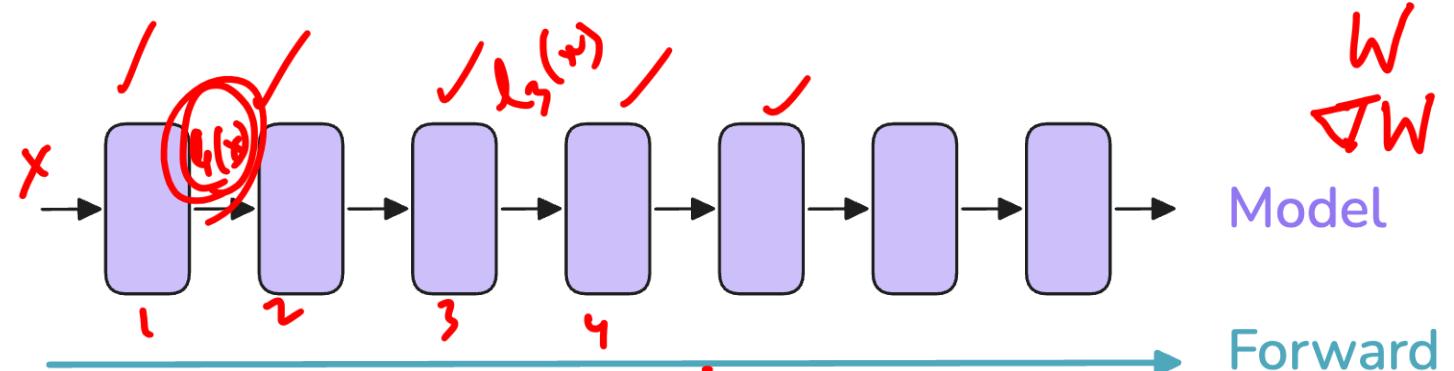


*Over 4,000 scaling experiments on up to 512 GPUs and measured throughput (size of markers) and GPU utilization (color of markers) - Normalized per model size in this visualization.*



$$y = \frac{Wx}{\|x\|_2} \quad \nabla W = \left( \frac{\partial L}{\partial y} \right)_{L \times 1} \frac{x^T}{\|x\|_2} dL \frac{\partial y}{\partial y}$$

# First Steps: Training on One GPU



- ✓ ① Model params.
- ✓ ② activations
- ✗ ③ Gradients
- ✓ ④ optim state (2)

(I) MLP (nn. modules)

$f_1 = \text{nn. Linear}$

$f_2 = \text{ReLU}$

$f_3 = \text{nn. Linear}'$

1. Forward Pass: pass inputs through the model to yield its outputs and compute Loss

forward(self, x)

for batch in dataloader:

batch = batch. cuda()

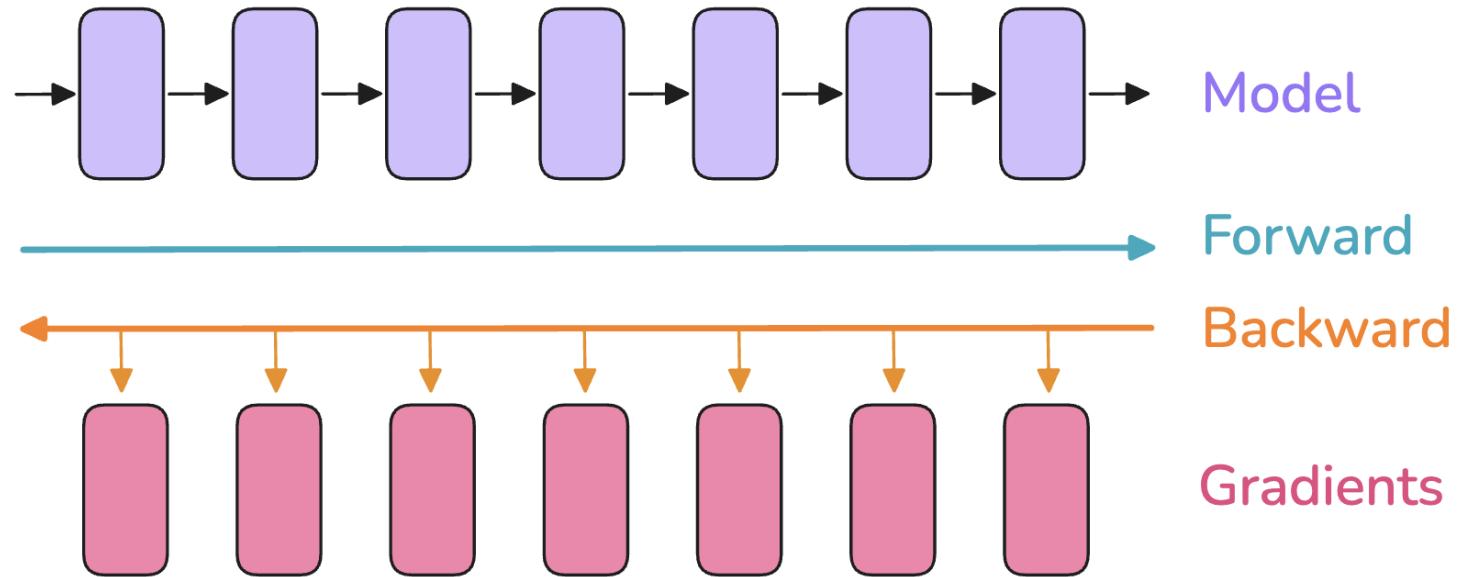
→ loss = model(batch)

→ loss.backward()

optim.step()

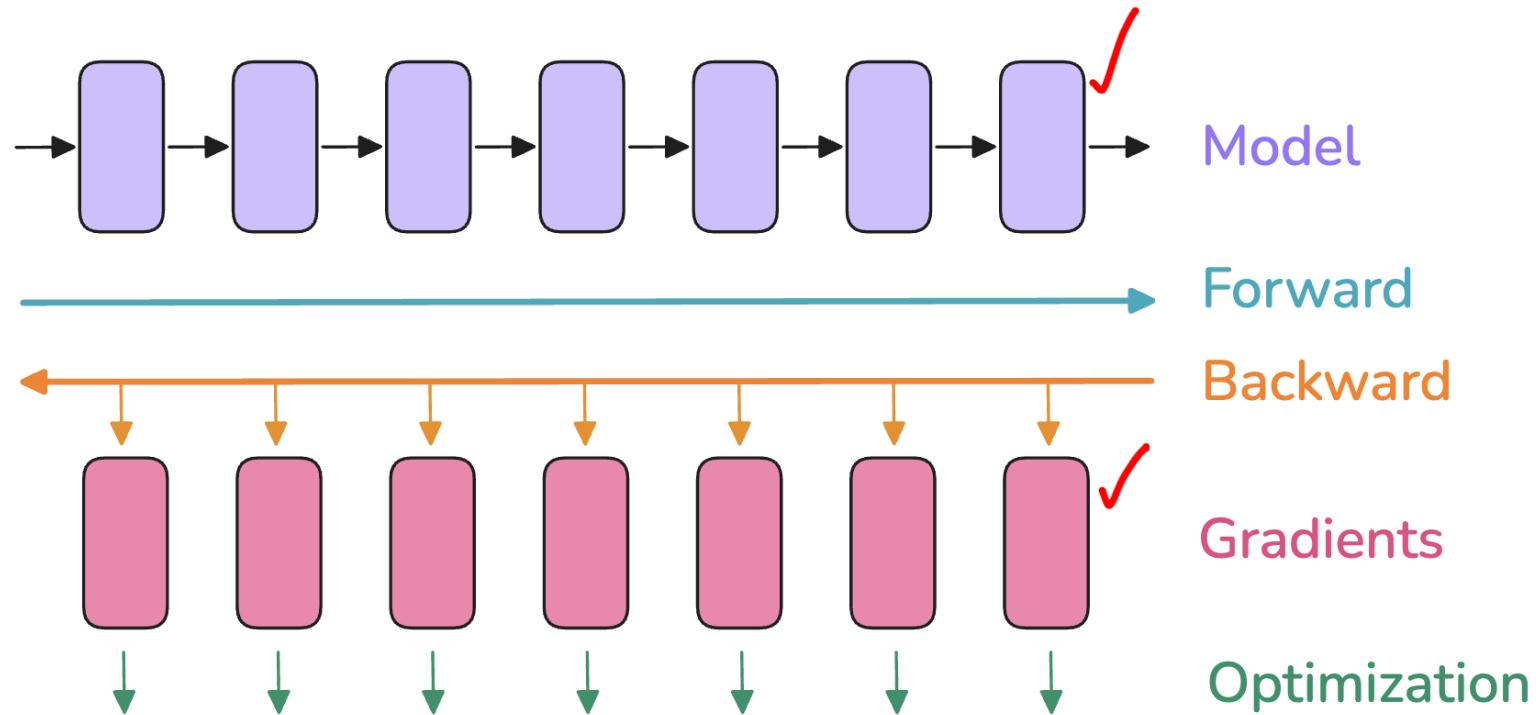


# First Steps: Training on One GPU



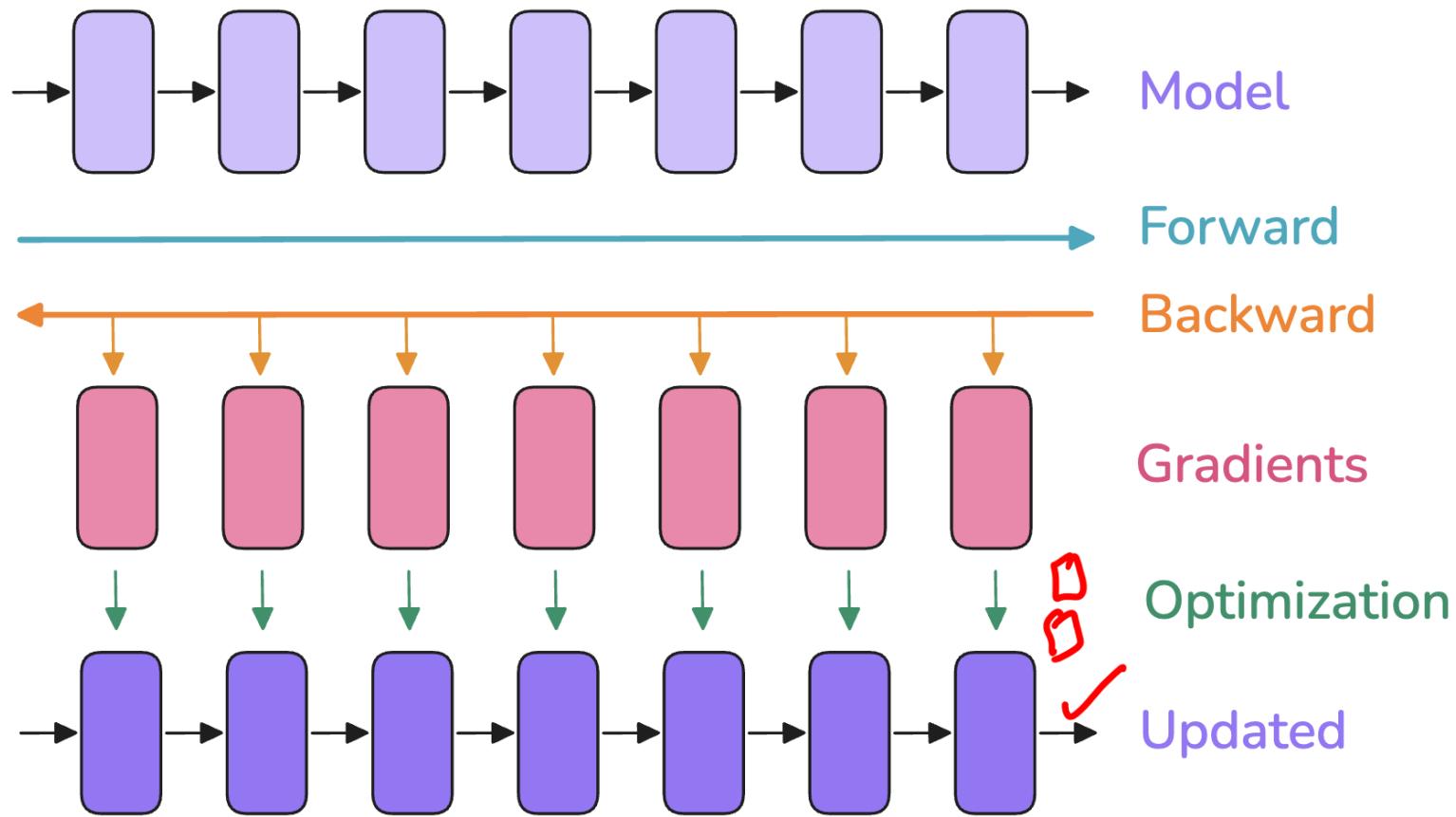
- 1. Forward Pass:** pass inputs through the model to yield its outputs and compute Loss
- 2. Backward Pass:** Compute the gradients

# First Steps: Training on One GPU



- 1. Forward Pass:** pass inputs through the model to yield its outputs and compute Loss
- 2. Backward Pass:** Compute the gradients
- 3. Optimization Step:** Update the parameters using the gradients

# First Steps: Training on One GPU

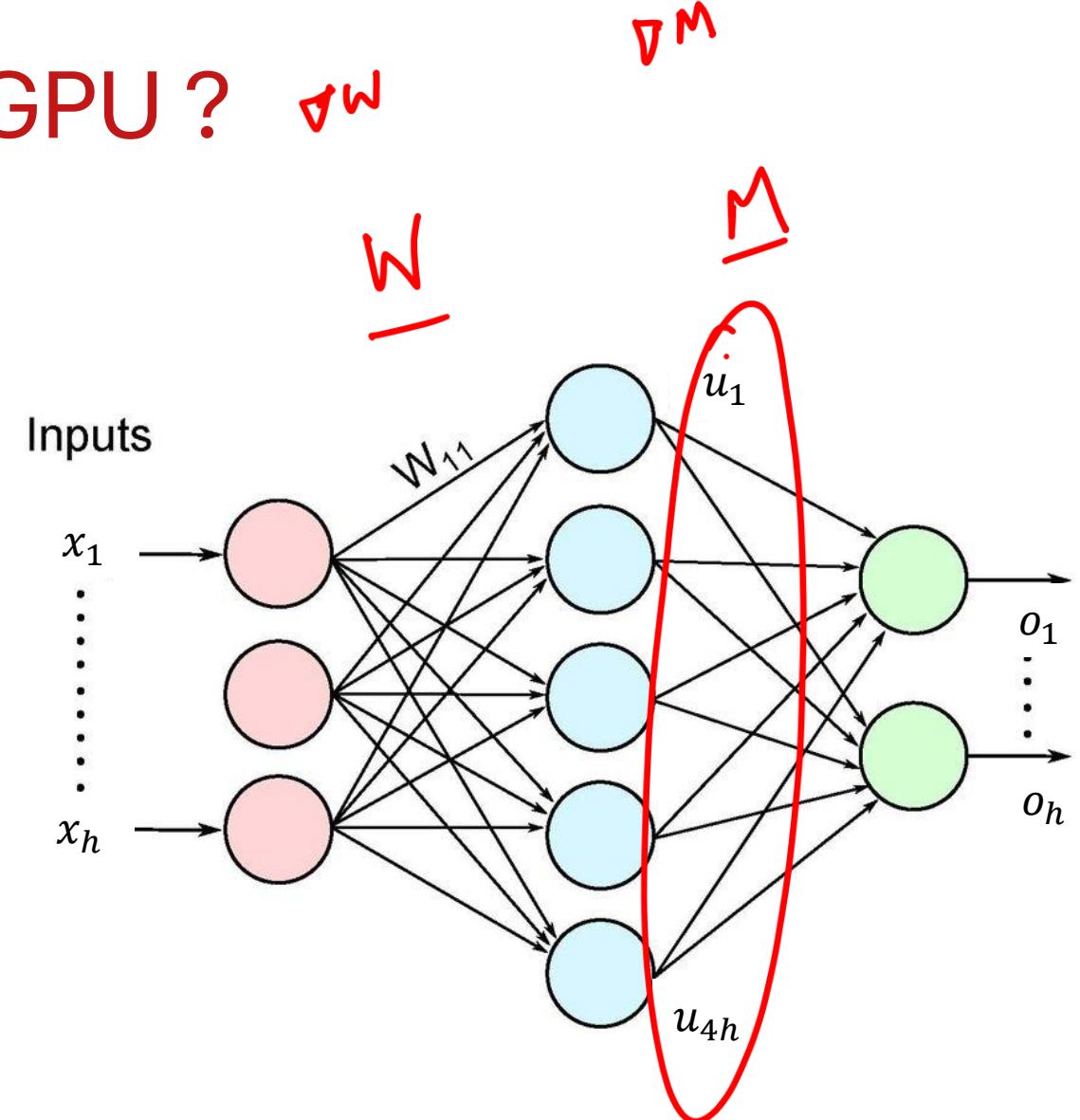


- 1. Forward Pass:** pass inputs through the model to yield its outputs and compute Loss
- 2. Backward Pass:** Compute the gradients
- 3. Optimization Step:** Update the parameters using the gradients

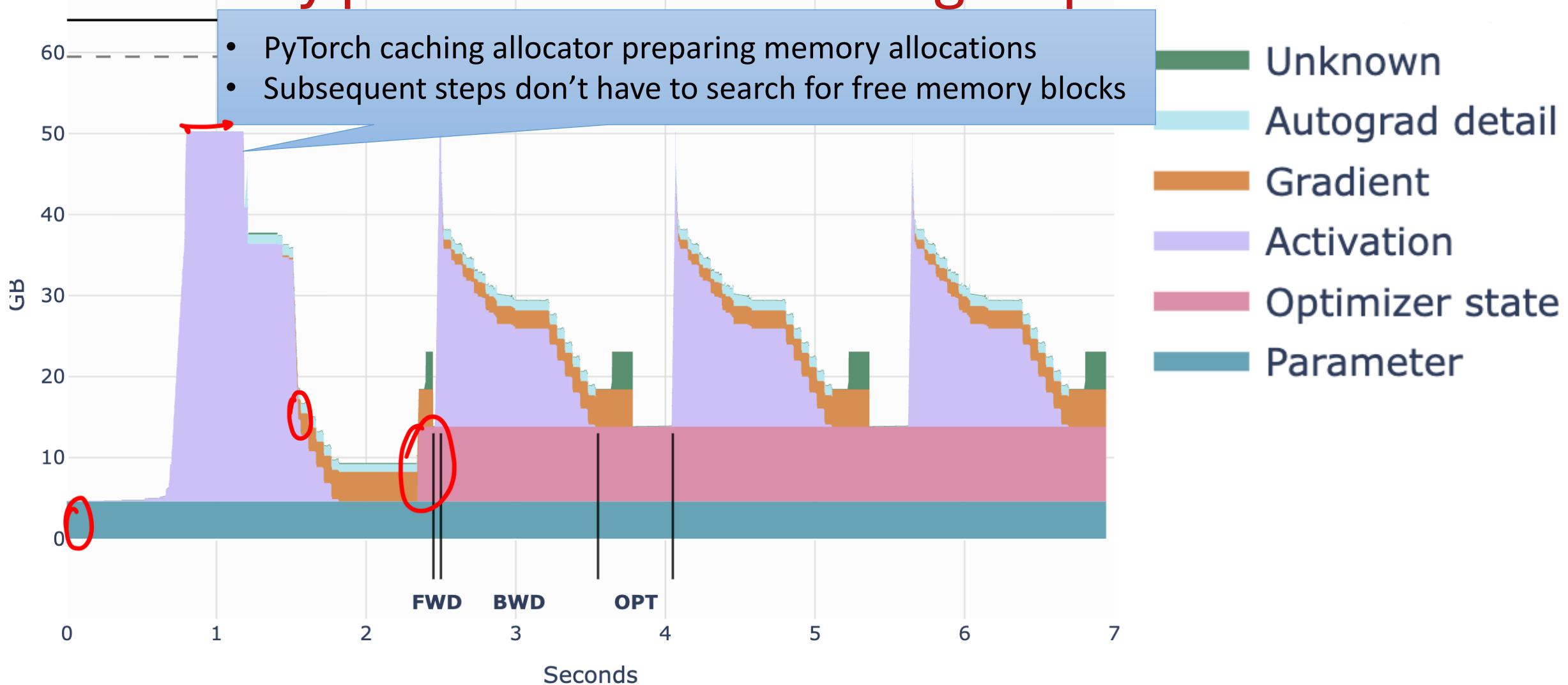
# First Steps: Training on One GPU ? $\sigma^W$

- **Memory usage in transformers**

- Model weights
- Model gradients
- Optimizer states
- Activations needed to compute the gradients

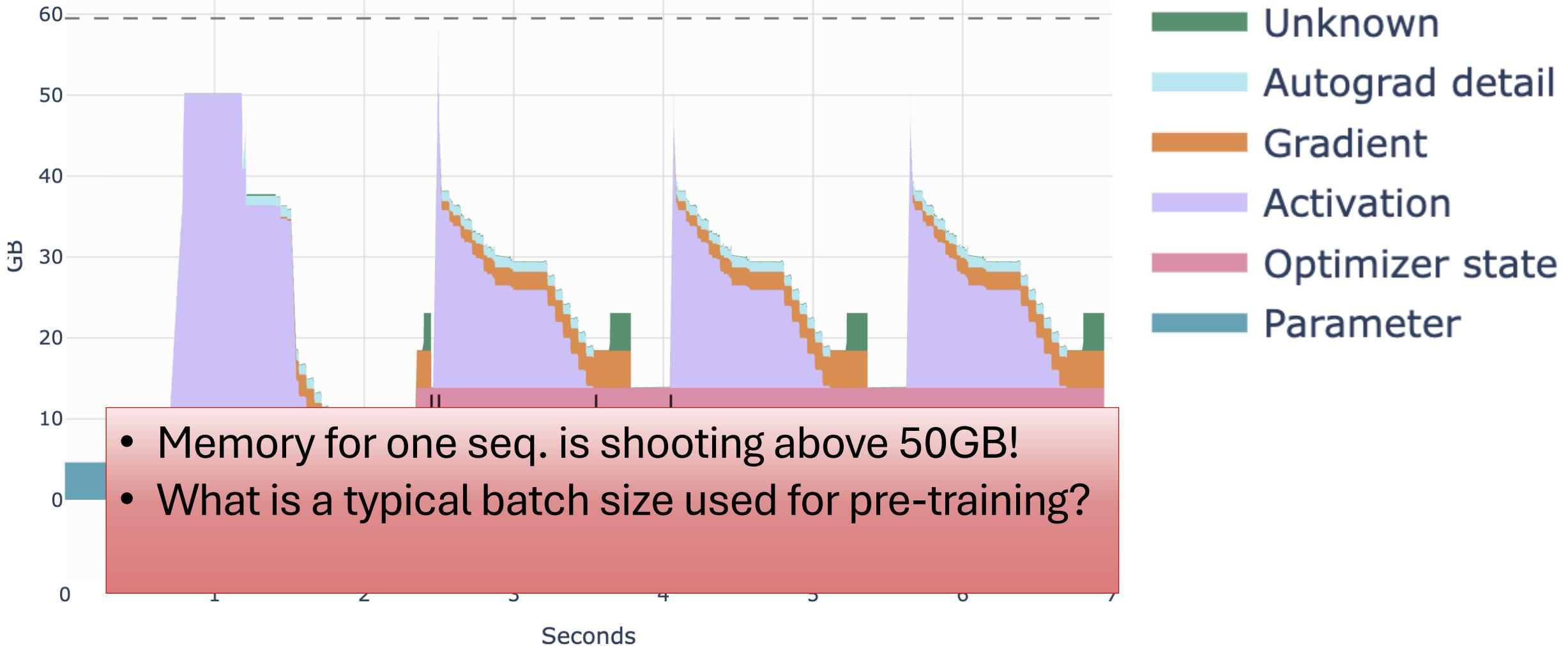


# Memory profile of first 4 training steps of Llama 1B



$$\text{th} = \frac{\text{seq.} \times \text{mbs}}{\text{lora} \times 400} \times 1024 \times 4$$

# Memory profile of first 4 training steps of Llama 1B



# Batch Size

- Reported in **#Tokens** = Batch Size Tokens =  $bst = bs * seq$   
*Llama1 – batch size ~ 4M for 1.4T tokens; DeepSeek - batch size ~ 60M for 14T tokens*
- Affects both model convergence and throughput

## Small batch size:

- 👎 Noisy gradient estimates
- 👍 Quick to compute => each step is fast
- 👎 More optimizer steps for fixed data
- 👎 May evade optima due to noisy gradients

## Large batch size:

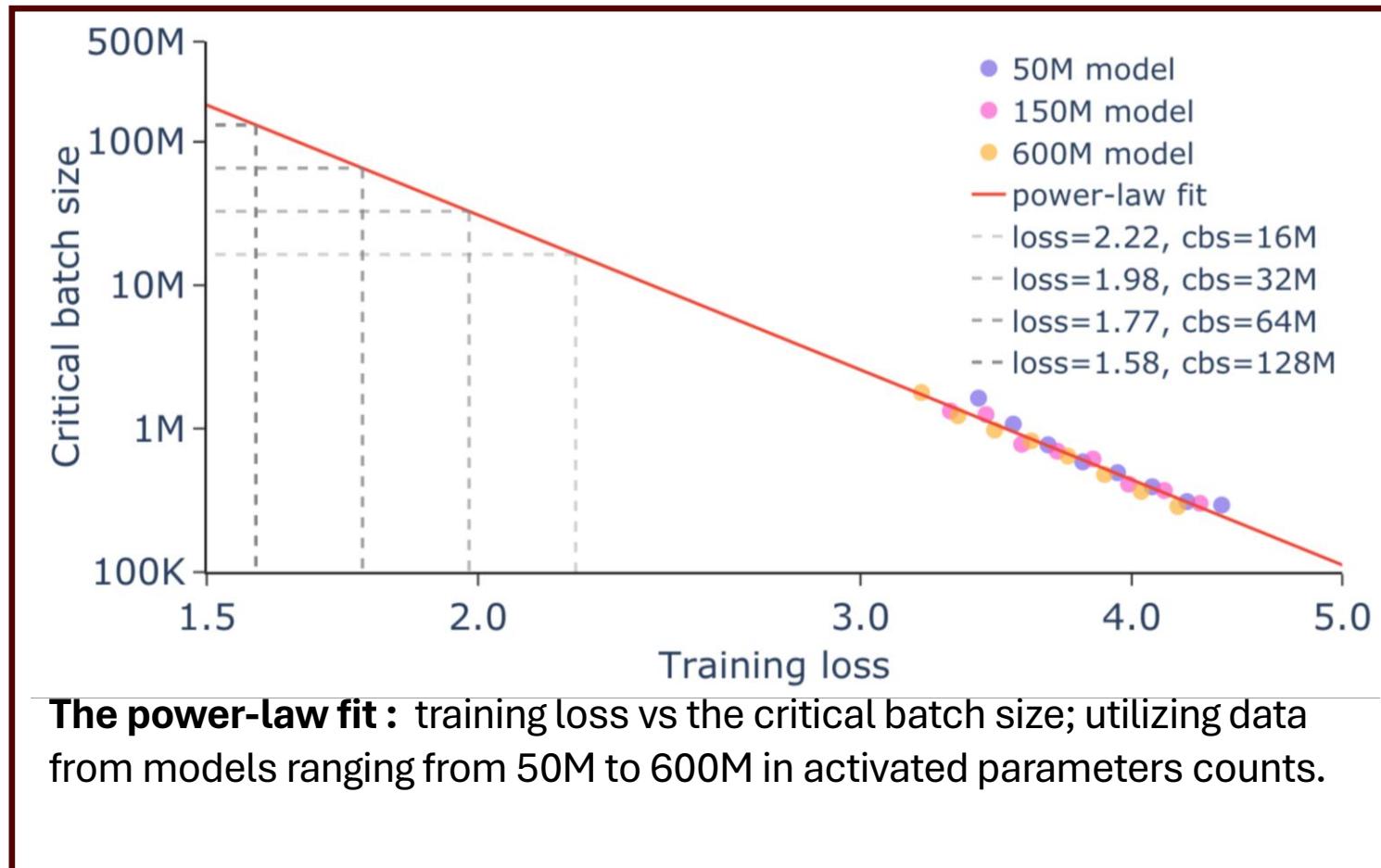
- 👍 Stable gradient estimates
- 👎 More time to compute => each step is slow
- 👍 Less optimizer steps for fixed data
- 👍 Arrive at better optima due to stable gradients



# Batch Size

- Training at the **critical batch size** yields a near-optimal balance between training time and data efficiency

[McCandlish et al. 2018, An Empirical Model of Large-Batch Training](#)



# Batch Size

## *Towards the beginning...*



### **Small Batch Size**

- Quickly move through the loss landscape  
(multiple noisy gradient update steps)



### **Large Batch Size**

- Accurate gradient, but slower convergence,  
potentially wasting compute



### **Small Batch Size**

Model may not converge to the most optimal performance (due to noisy gradients)



### **Large Batch Size**

Accurate gradient, helps in reaching optimal performance

**Batch Size:** 16M

32M

64M

128M

**# Tokens Seen:** 0

69B

790B

4700B

10,400B

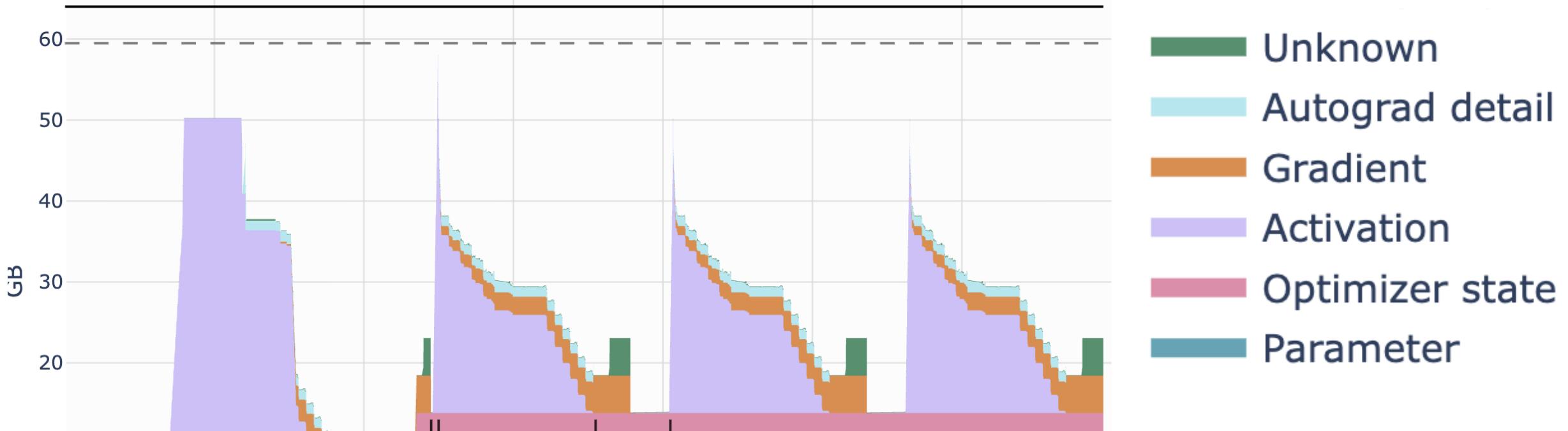
**Gradient update steps →**

[https://filecdn.minimax.chat/\\_Arxiv\\_Minimax\\_01\\_Report.pdf](https://filecdn.minimax.chat/_Arxiv_Minimax_01_Report.pdf)

Kaplan et al. 2020 Scaling laws for neural language models.



# Memory profile of first 4 training steps of Llama 1B



- Memory for one seq. is shooting above 50GB!
- What is a typical batch size used for pre-training?
- How much memory a transformer layer uses?

$$x \leftarrow \mu + \sigma \left( \frac{x - \mathbb{E}(x)}{\text{Cov}(x)} \right)$$

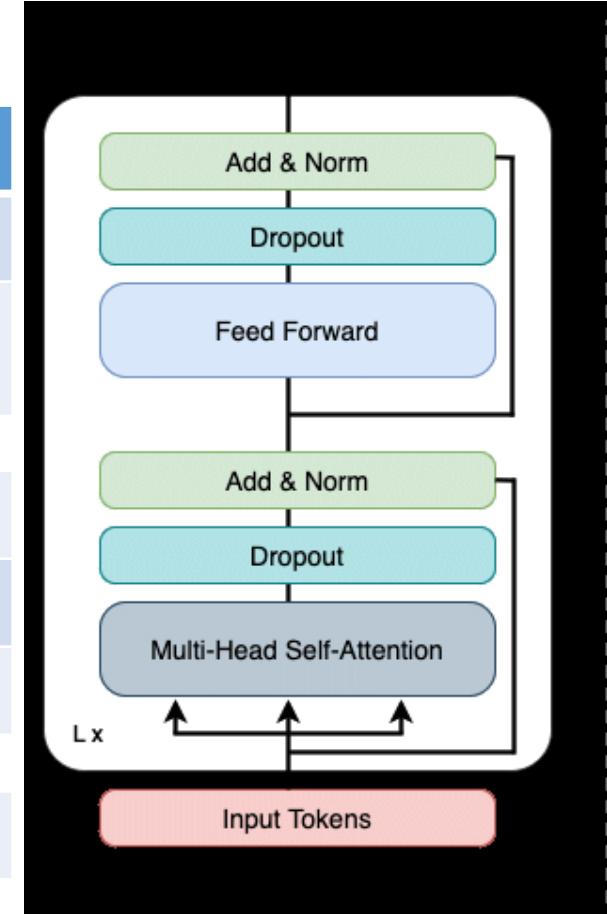
# Memory for Weights, Grads, and Optim States



- ✓  $h$  : hidden dimension,
- ✓  $v$  : vocabulary size,
- $L$  : number of layers.

Total	
Layer Norm	$2h$
MLP	$(h \times 4h + 4h) + 4h^2 + h$

Layer Norm	$2h$
Projection	$h^2 + h$
MHA	$3(h^2 + h)$
Token Emb.	$hv$



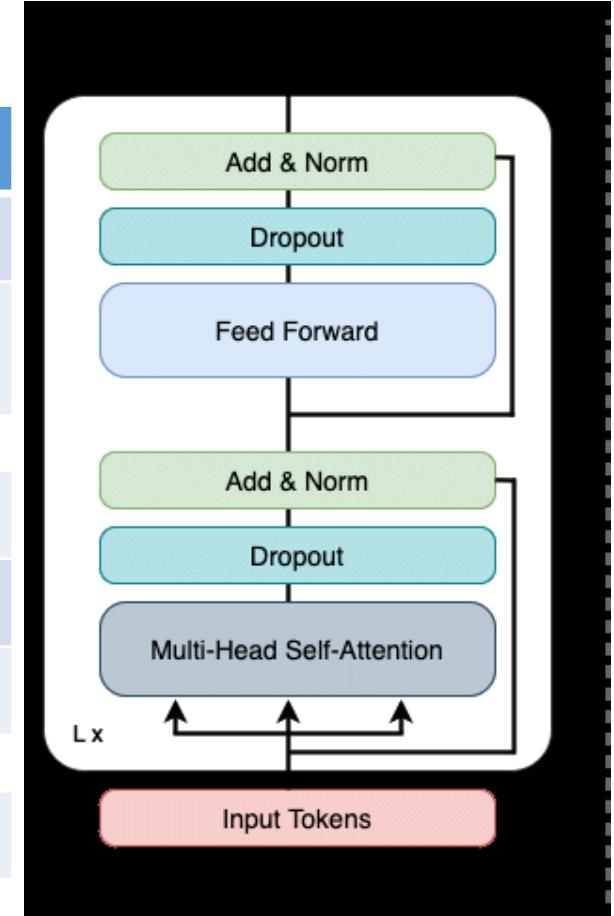
<https://michaelwornow.net/2024/01/18/counting-params-in-transformer>

# Memory for Weights, Grads, and Optim States

$$\bullet N = \underbrace{h * v}_{\text{Total}} + \underbrace{L * (12 * h^2 + 13 * h)}_{\text{MLP}} + \underbrace{2 * h}_{\text{Layer Norm}}$$

$h$  : hidden dimension,  
 $v$  : vocabulary size,  
 $L$  : number of layers.

Total	
Layer Norm	$2 * h$
MLP	$h * 4h + 4h + 4h * h + h$
Layer Norm	$2 * h$
Projection	$h * h + h$
MHA	$3 * (h * h + h)$
Token Emb.	$h * v$



<https://michaelwornow.net/2024/01/18/counting-params-in-transformer>



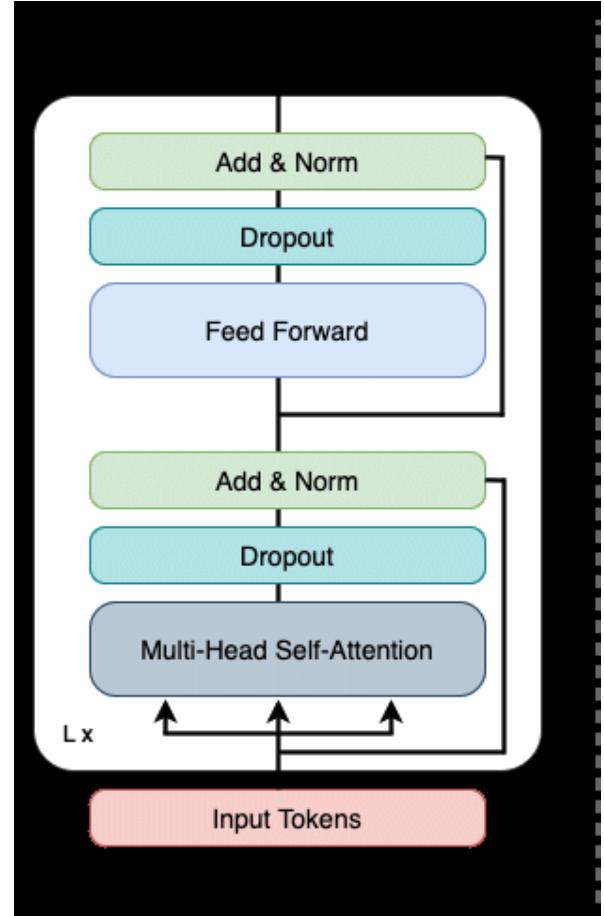
# Memory for Weights, Grads, and Optim States

- $N = h * v + L * (12 * h^2 + 13 * h) + 2 * h$
- FP32 representation ( 4 bytes per param)
  - $m_{params} = 4 * N$
  - $m_{grad} = 4 * N$
  - $m_{opt} = (4 + 4) * N$
- Mixed precision: computation in bf16 (2 bytes); storage in FP32
  - $m_{params} = 2 * N$
  - $m_{grad} = 2 * N$
  - $m_{opt} = (4 + 4) * N$
  - $m_{params\_fp32} = 4 * N$

Total:  $16 * N$  bytes

Total:  $16 * N$  bytes

Why use mixed precision if total memory is same?



# Memory for Weights, Grads, and Optim States

- $N = h * v + L * (12 * h^2 + 13 * h) + 2 * h$

- FP32 representation ( 4 bytes per param)

- $m_{params} = 4 * N$
- $m_{grad} = 4 * N$
- $m_{opt} = (4 + 4) * N$

Total:  $16 * N$  bytes

- Mixed precision: computation in bf16 (2 bytes); storage in FP32

- $m_{params} = 2 * N$
- $m_{grad} = 2 * N$
- $m_{opt} = (4 + 4) * N$
- $m_{params\_fp32} = 4 * N$   
(master weights)

Total:  $16 * N$  bytes

1. Allows us to use optimized lower precision operations on the GPU, which are faster
2. Reduces the activation memory requirements during the forward pass

Why use mixed precision if total memory is same?



# Memory for Weights, Grads, and Optim States

$$\bullet N = h * v + L * (12 * h^2 + 13 * h) + 2 * h$$

- FP32 representation ( 4 bytes per param)

- $m_{params} = 4 * N$
- $m_{grad} = 4 * N$
- $m_{opt} = (4 + 4) * N$

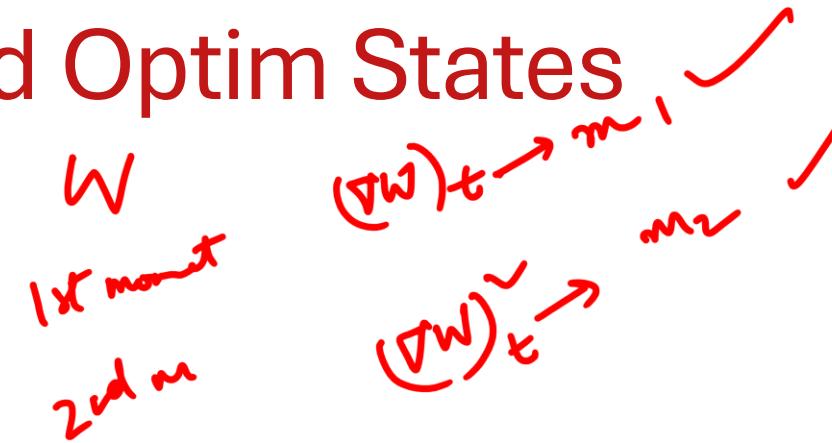
Total:  $16 * N$  bytes

- Mixed precision: computation in bf16 (2 bytes); storage in FP32

- $m_{params} = 2 * N$
- $m_{grad} = 2 * N$
- $m_{opt} = (4 + 4) * N$
- $m_{params\_fp32} = 4 * N$   
(master weights)

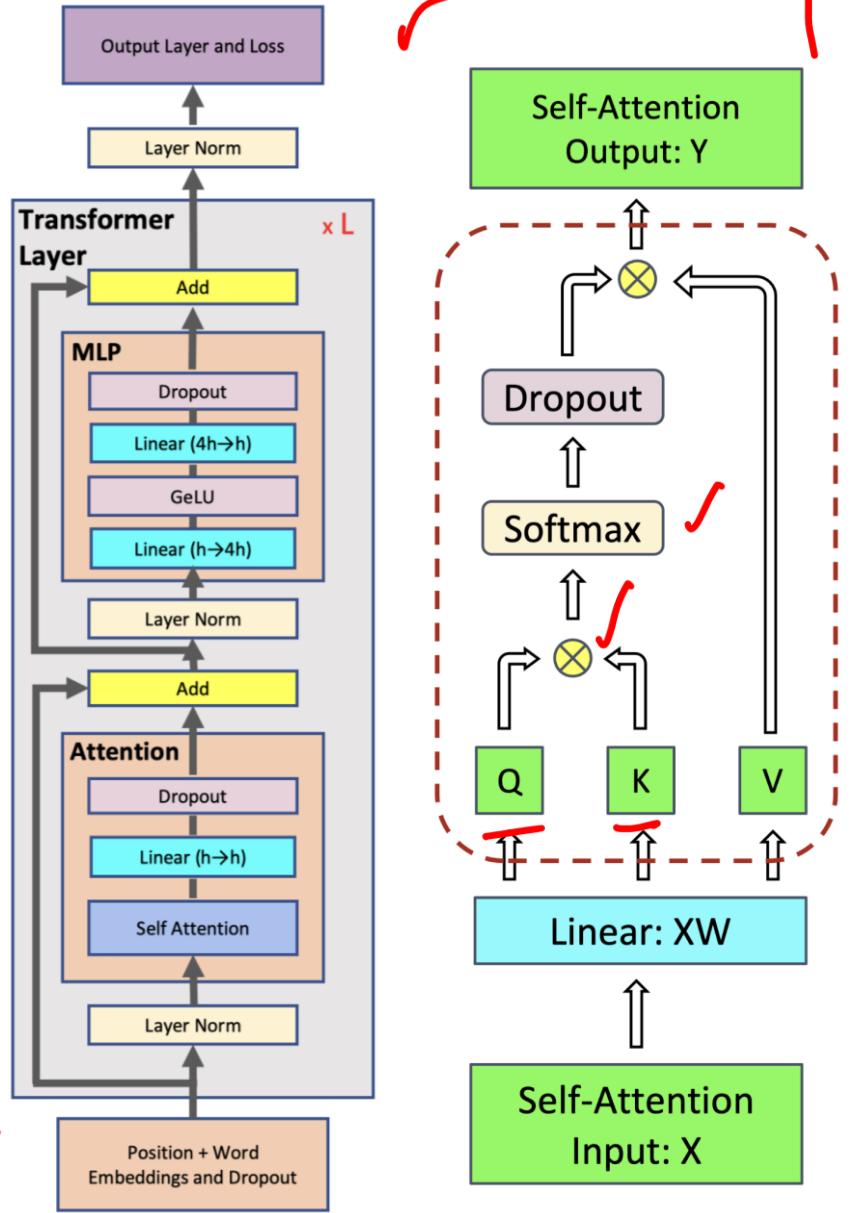
Total:  $16 * N$  bytes

Why use mixed precision if total memory is same?



Model parameters	Total Mem.
1B	16 GB
7B	112 GB
70B	1120 GB
405B	6480 GB

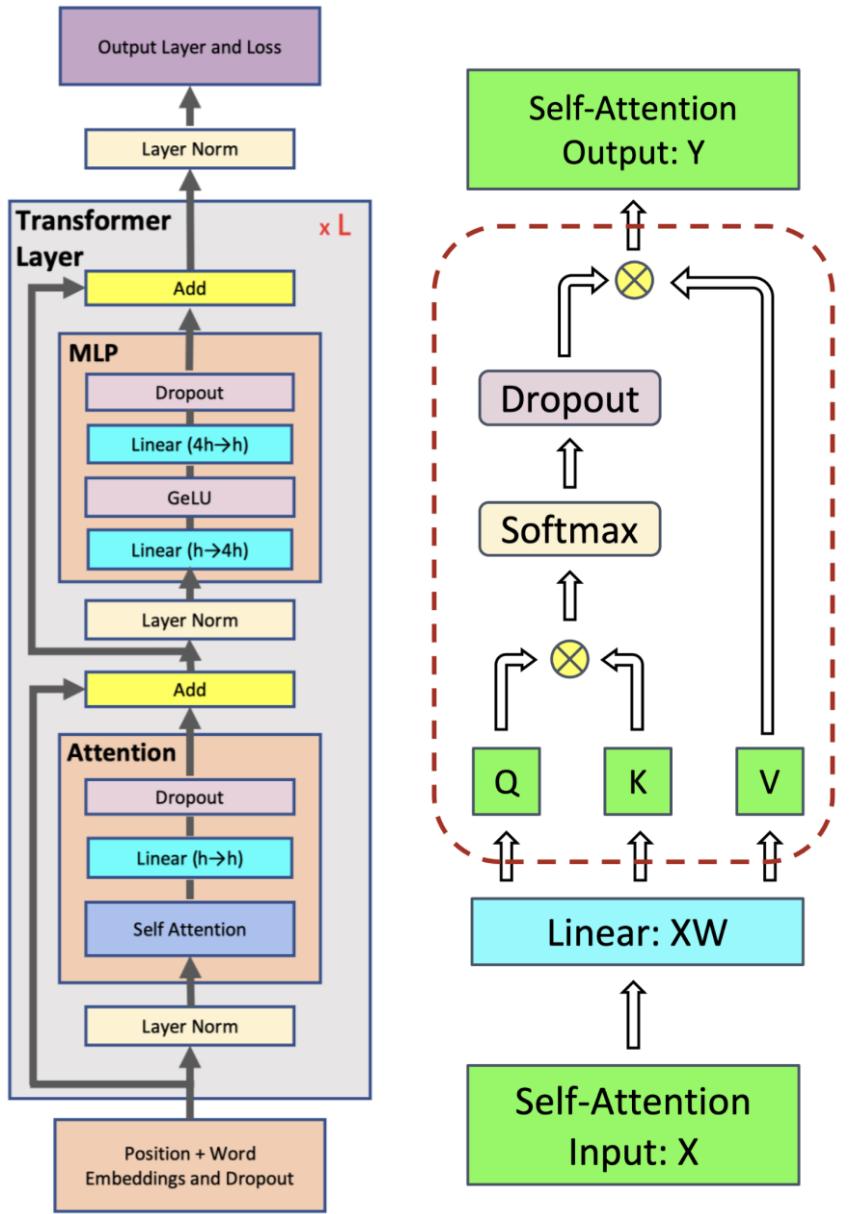




# Memory for Activations

MHA	
$Q, K, V$ input	$bs \times seq \times h$
$QK^T$	$O(b \cdot sh)$
Softmax	$O(n_{heads} \cdot seq \cdot bs)$
Softmax d/o mask	
Prob*Values	

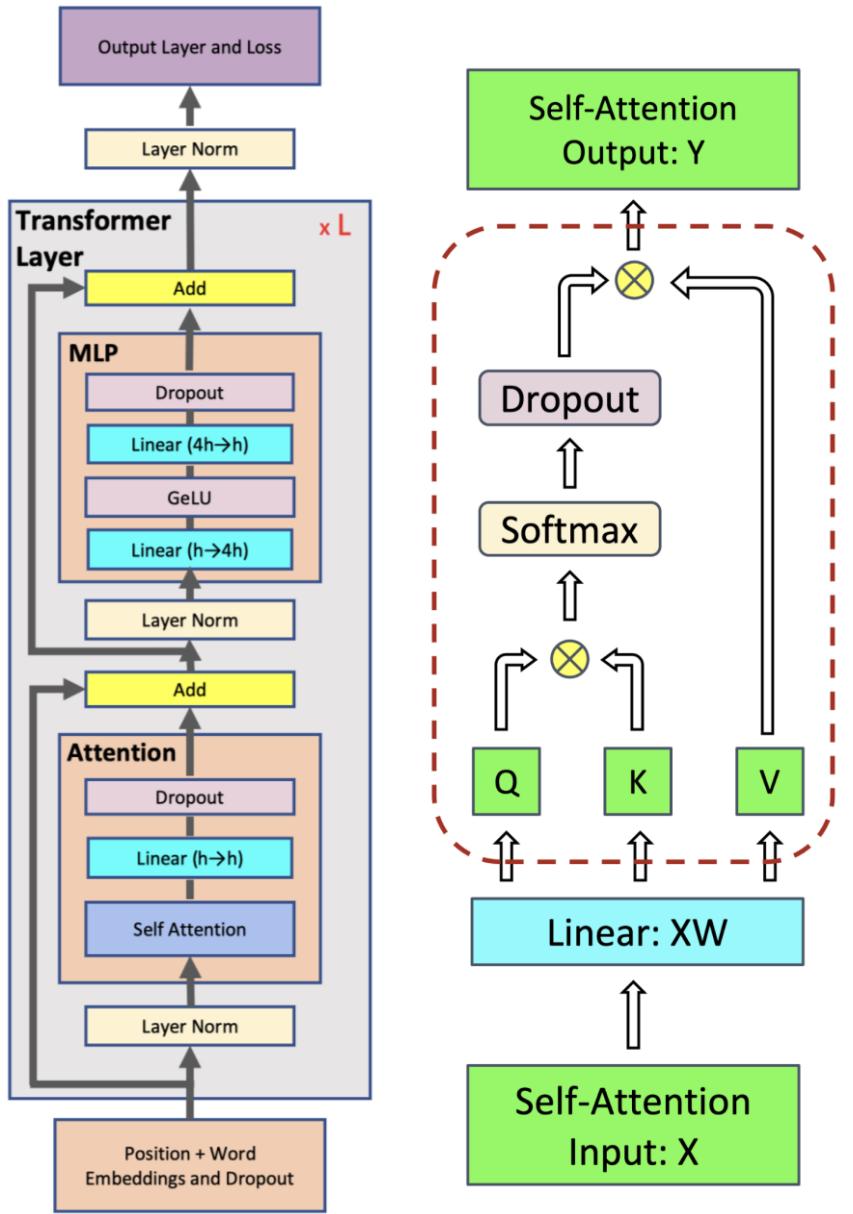
Korthikanti et al. 2022, Reducing Activation Recomputation in Large Transformer Models



# Memory for Activations

Attention Block	
<b>Dropout mask</b>	
<b>Linear Proj:</b>	
<b>MHA</b>	$8 * \text{seq} * \text{bs} * h + 5 * \text{nheads} * \text{seq}^2 * \text{bs}$
<b><math>Q, K, V</math> input</b>	$2 * \text{seq} * \text{bs} * h$
<b><math>QK^T</math></b>	$4 * \text{seq} * \text{bs} * h$
<b>Softmax</b>	$2 * \text{nheads} * \text{bs} * \text{seq} * \text{seq}$
<b>Softmax d/o mask</b>	$\text{nheads} * \text{bs} * \text{seq} * \text{seq}$
<b>Prob*Values</b>	$2 * \text{nheads} * \text{bs} * \text{seq} * \text{seq} + 2 * \text{seq} * \text{bs} * h$

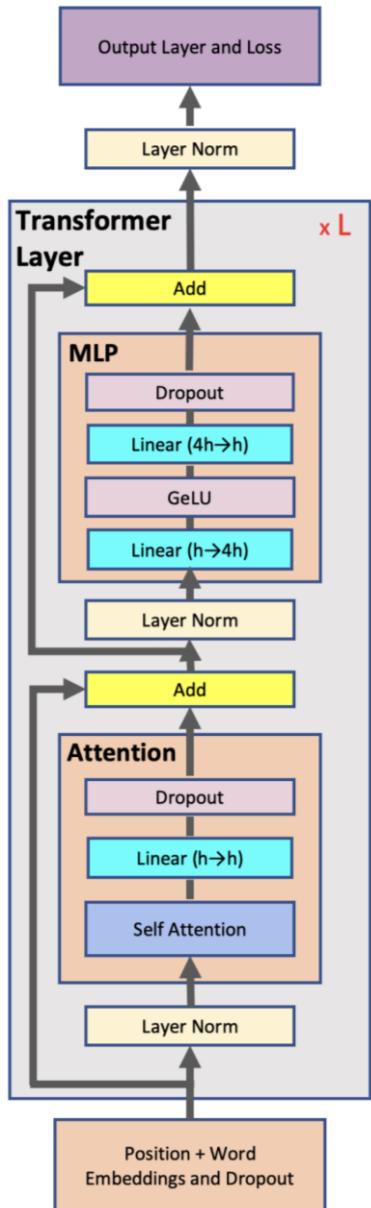
Korthikanti et al. 2022, Reducing Activation Recomputation in Large Transformer Models



# Memory for Activations

<b>Attention Block</b>	$11 * seq * bs * h + 5 * n_{heads} * seq^2 * bs$
<b>Dropout mask</b>	$seq * bs * h$
<b>Linear Proj:</b>	$2 * seq * bs * h$
<b>MHA</b>	$8 * seq * bs * h + 5 * n_{heads} * seq^2 * bs$
<b><math>Q, K, V</math> input</b>	$2 * seq * bs * h$
<b><math>QK^T</math></b>	$4 * seq * bs * h$
<b>Softmax</b>	$2 * n_{heads} * bs * seq * seq$
<b>Softmax d/o mask</b>	$n_{heads} * bs * seq * seq$
<b>Prob*Values</b>	$2 * n_{heads} * bs * seq * seq + 2 * seq * bs * h$

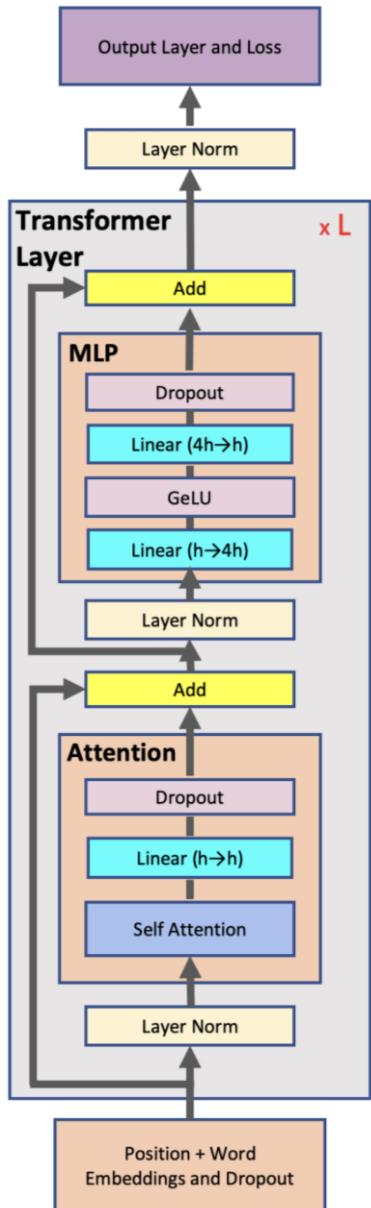
Korthikanti et al. 2022, Reducing Activation Recomputation in Large Transformer Models



# Memory for Activations

<b>MLP Block</b>	
D/o mask	
Linear ( $4h \rightarrow h$ )	
GELU	
Linear ( $h \rightarrow 4h$ )	
<b>Layer Norm</b>	
<b>Attention Block</b>	$11 * seq * bs * h +$ $5 * n_{heads} * seq^2 * bs$
<b>Layer Norm</b>	

Korthikanti et al. 2022, Reducing Activation Recomputation in Large Transformer Models



<b>Total</b>	$34 * \text{seq} * \text{bs} * h + 5 * n_{\text{heads}} * \text{seq}^2 * \text{bs}$
<b>MLP Block</b>	$19 * \text{seq} * \text{bs} * h$
D/o mask	$1 * \text{seq} * \text{bs} * h$
Linear ( $4h \rightarrow h$ )	$8 * \text{seq} * \text{bs} * h$
GELU	$8 * \text{seq} * \text{bs} * h$
Linear ( $h \rightarrow 4h$ )	$2 * \text{seq} * \text{bs} * h$
<b>Layer Norm</b>	$2 * \text{seq} * \text{bs} * h$
<b>Attention Block</b>	$11 * \text{seq} * \text{bs} * h + 5 * n_{\text{heads}} * \text{seq}^2 * \text{bs}$
<b>Layer Norm</b>	$2 * \text{seq} * \text{bs} * h$

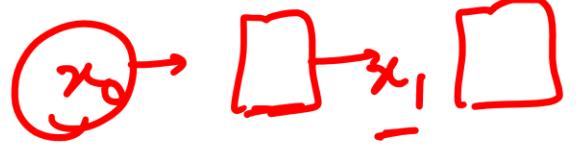
# Memory for Activations

$m_{\text{act}}$

$$= L * \left( 34 * \text{seq} * \text{bs} * h + 5 * n_{\text{heads}} * \text{seq}^2 * \text{bs} \right)$$

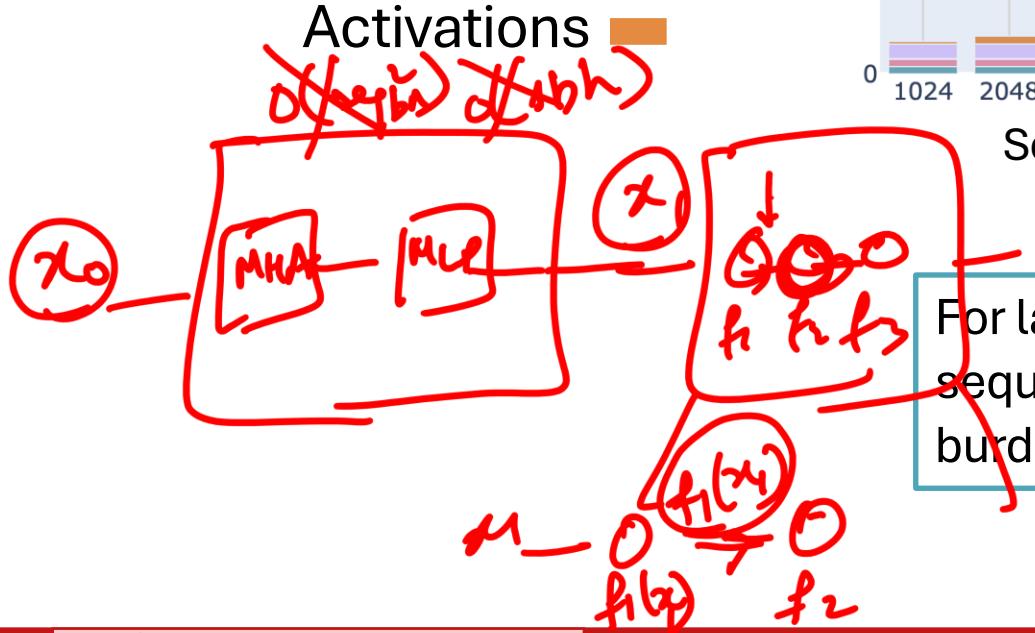
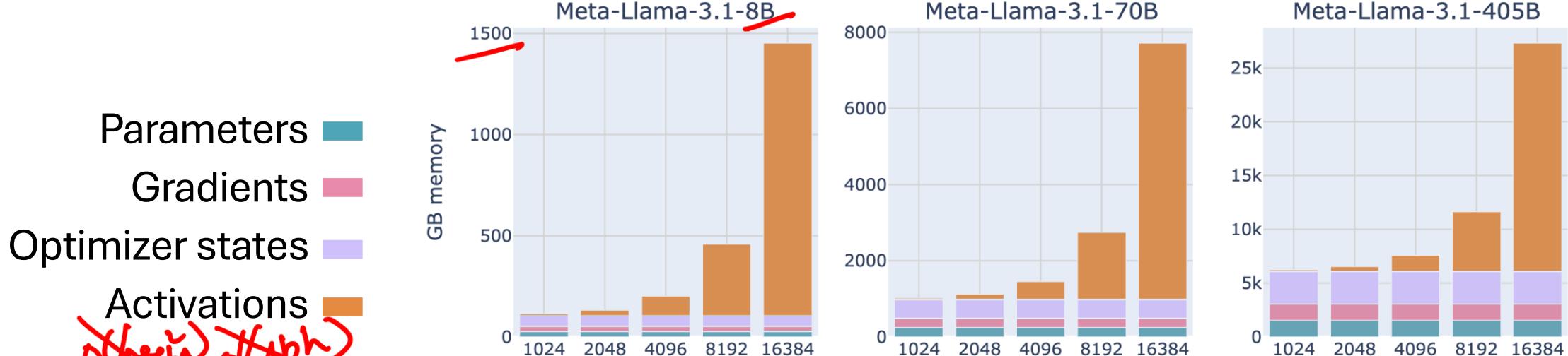
- Scales linearly with batch size
- Quadratically with the sequence length

$$y = Wx$$



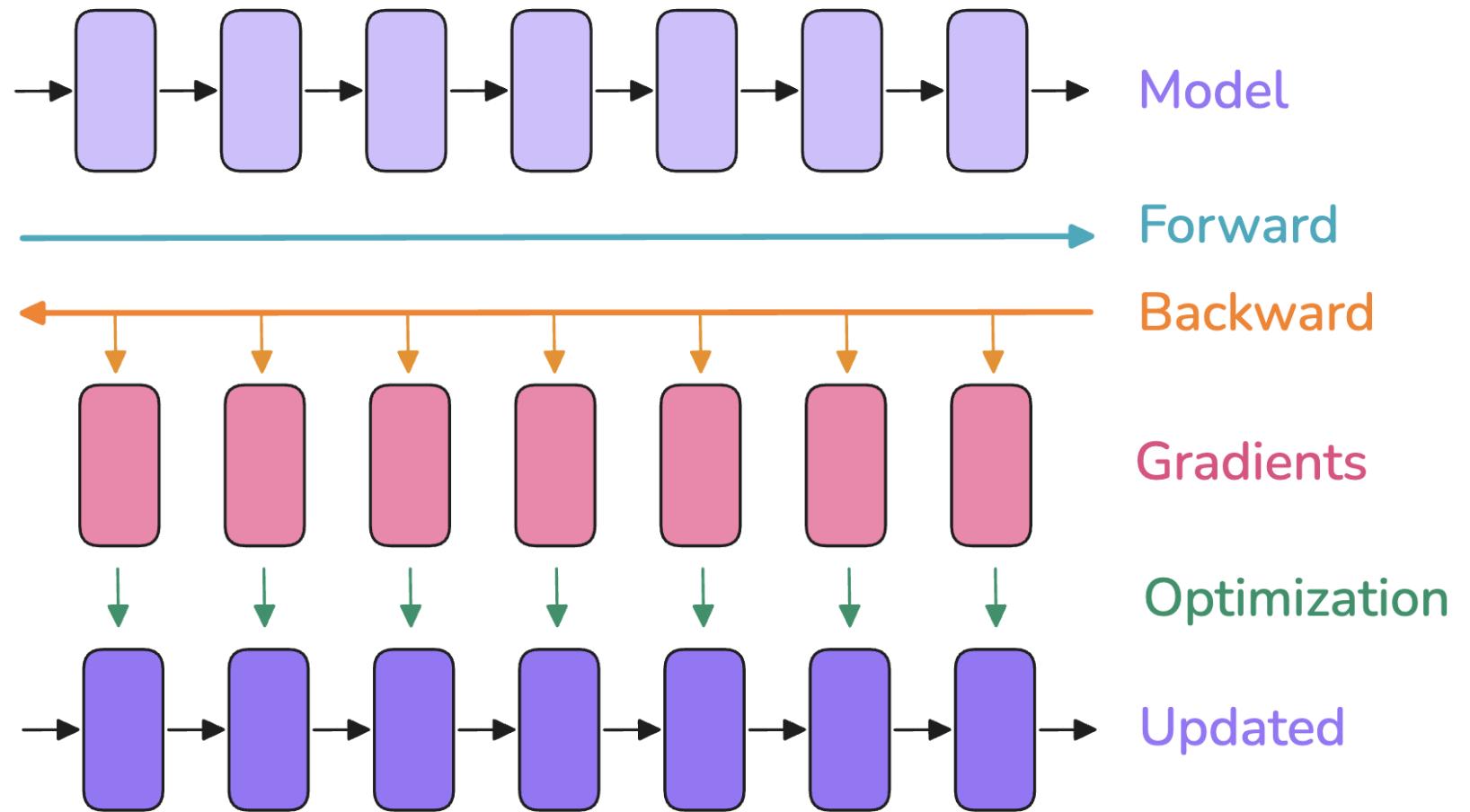
# Memory usage in Transformers

Batch Size: 1



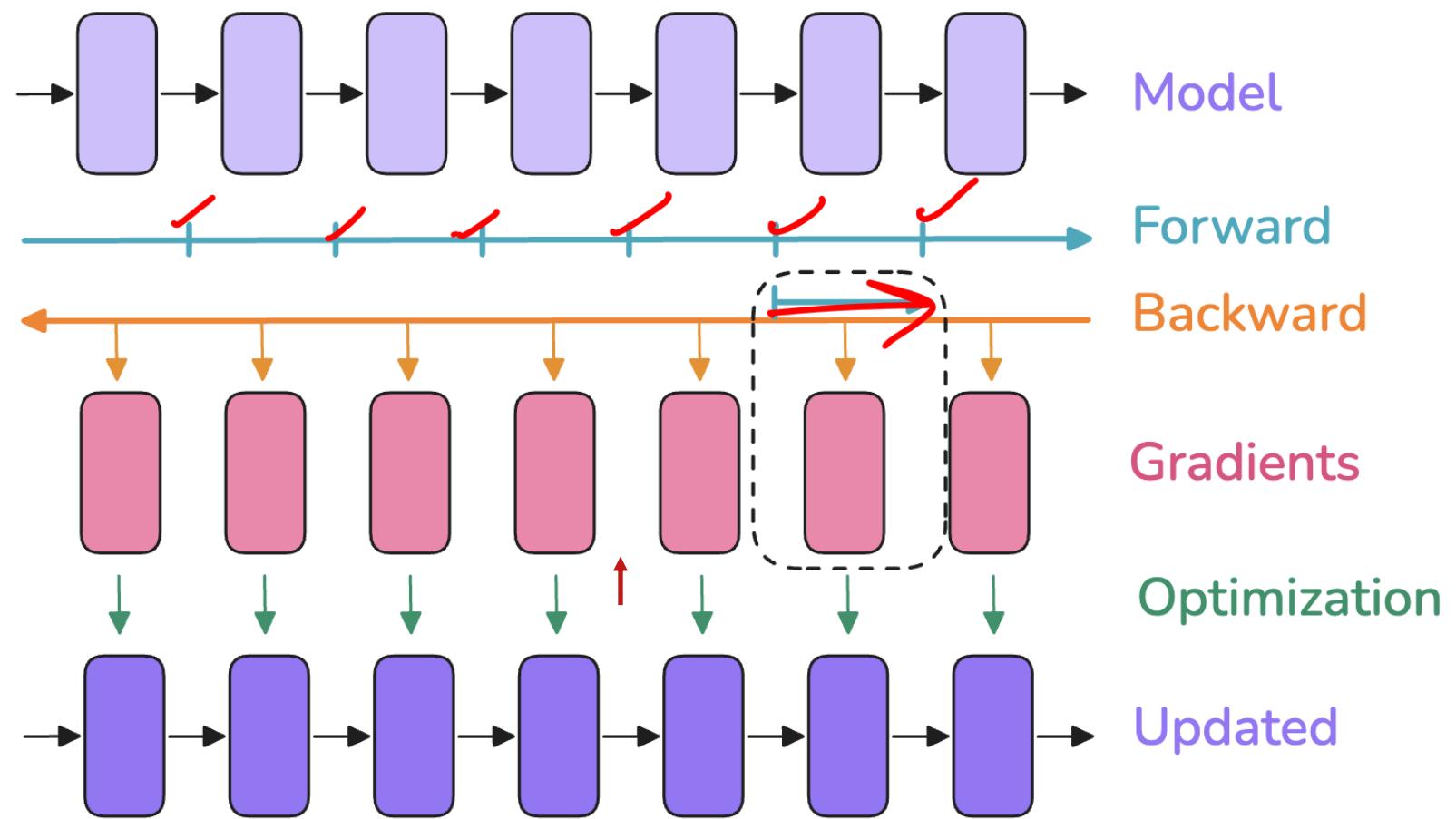
For large numbers of input tokens (i.e., large batch sizes / sequences), activations become by far the largest memory burden.

# Activation re-computation aka Gradient Checkpointing



# Activation re-computation aka Gradient Checkpointing

- Discard some activations during the forward pass to save memory
- Spend some extra compute to recompute these on the fly during the backward pass.



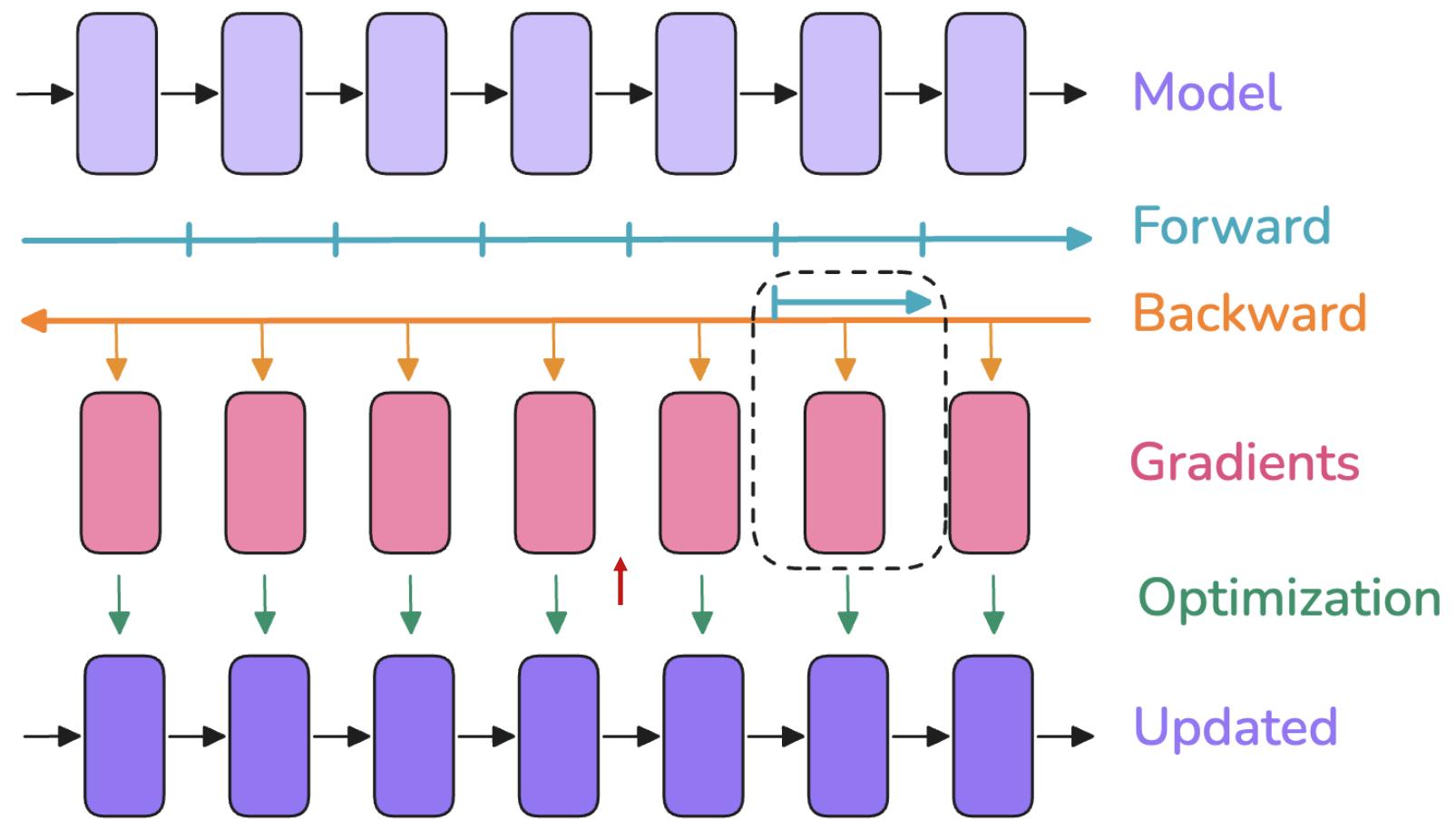
[Korthikanti et al. 2022, Reducing Activation Recomputation in Large Transformer Models](#)



# Activation re-computation aka Gradient Checkpointing

- Discard some activations during the forward pass to save memory
- Spend some extra compute to recompute these on the fly during the backward pass.

Which activations to store?



Korthikanti et al. 2022, Reducing Activation Recomputation in Large Transformer Models

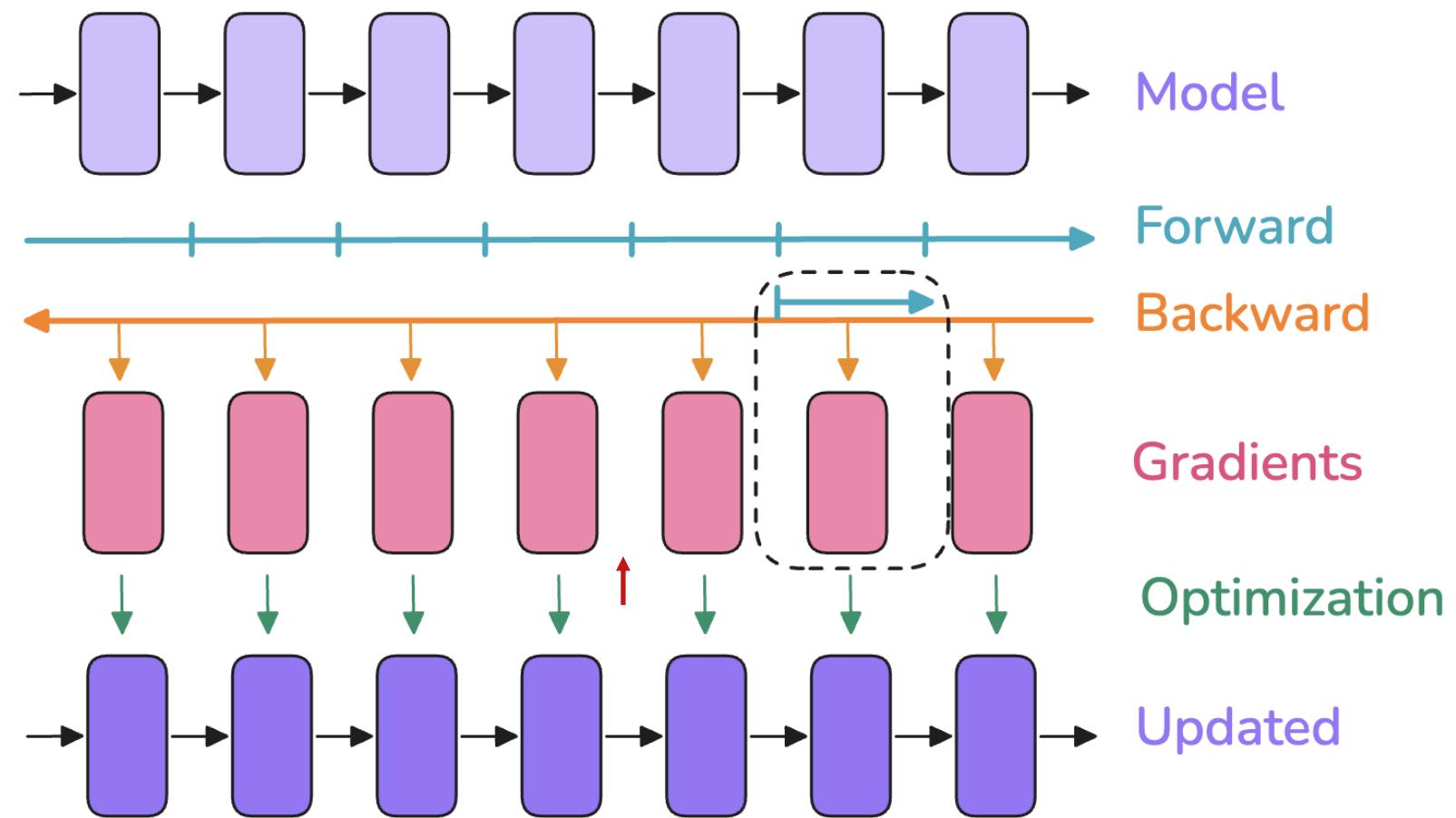


# Activation re-computation aka Gradient Checkpointing

- Discard some activations during the forward pass to save memory
- Spend some extra compute to recompute these on the fly during the backward pass.

## Which activations to store?

- **Full:** Save at the transition of layers
  - 👍 Saves most memory
  - 👎 Most expensive – 30-40%



[Korthikanti et al. 2022, Reducing Activation Recomputation in Large Transformer Models](#)

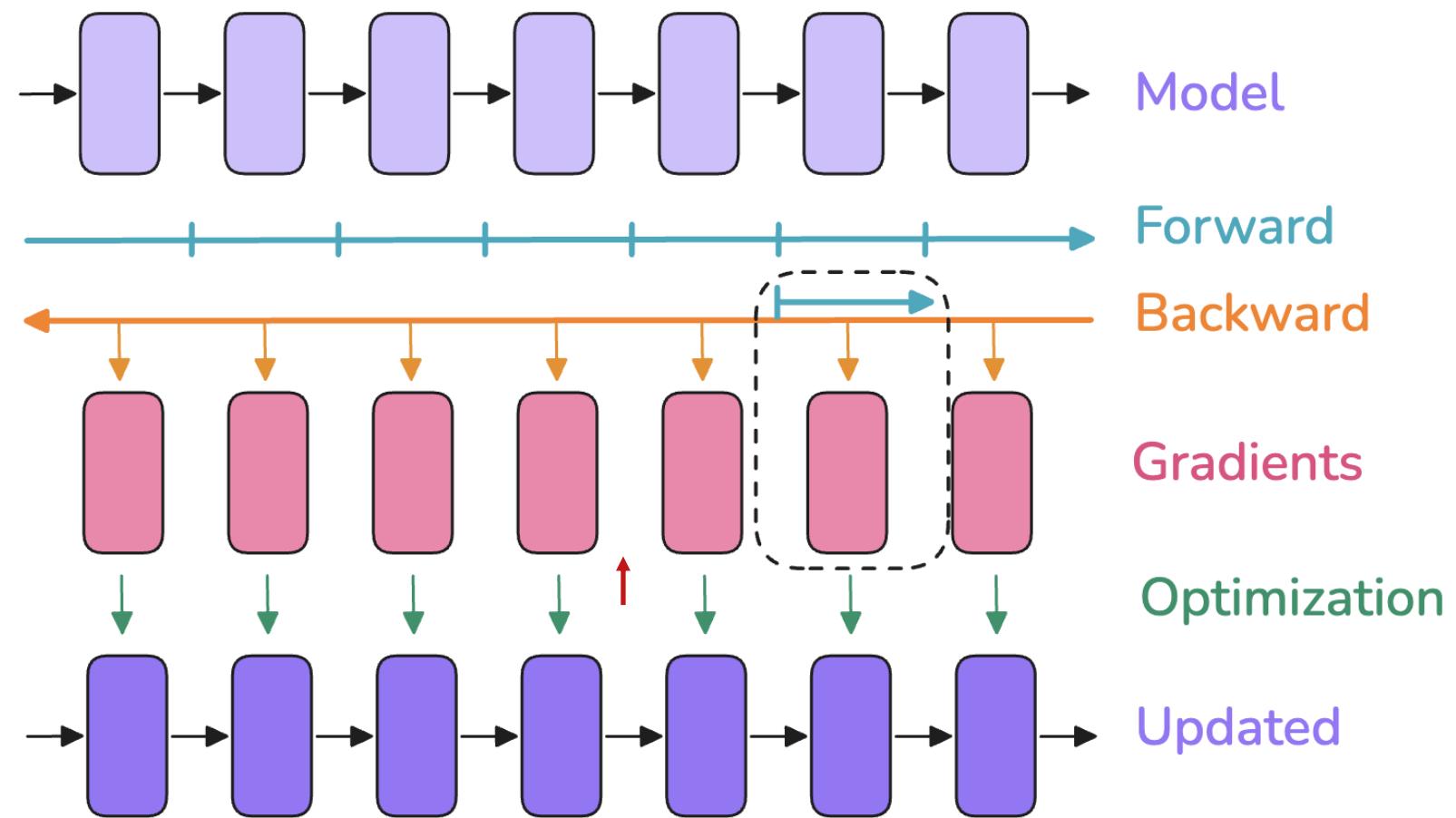


# Activation re-computation aka Gradient Checkpointing

- Discard some activations during the forward pass to save memory
- Spend some extra compute to recompute these on the fly during the backward pass.

## Which activations to store?

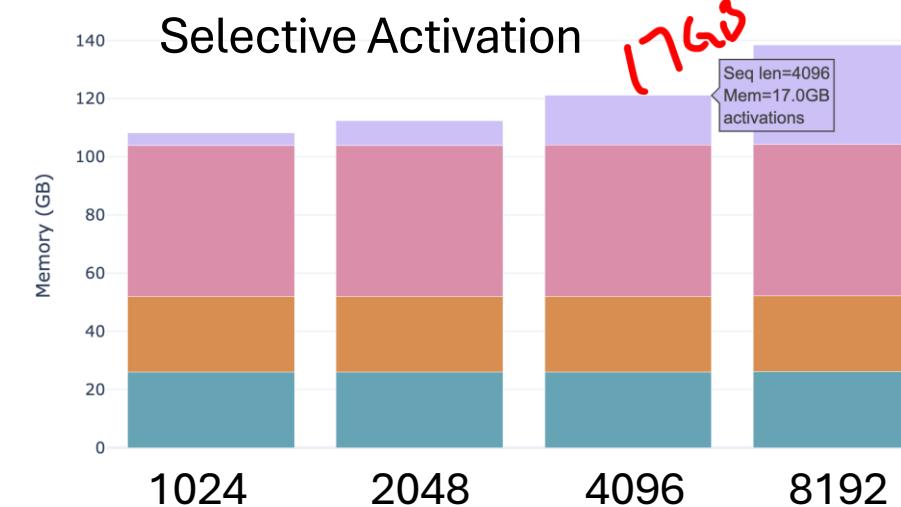
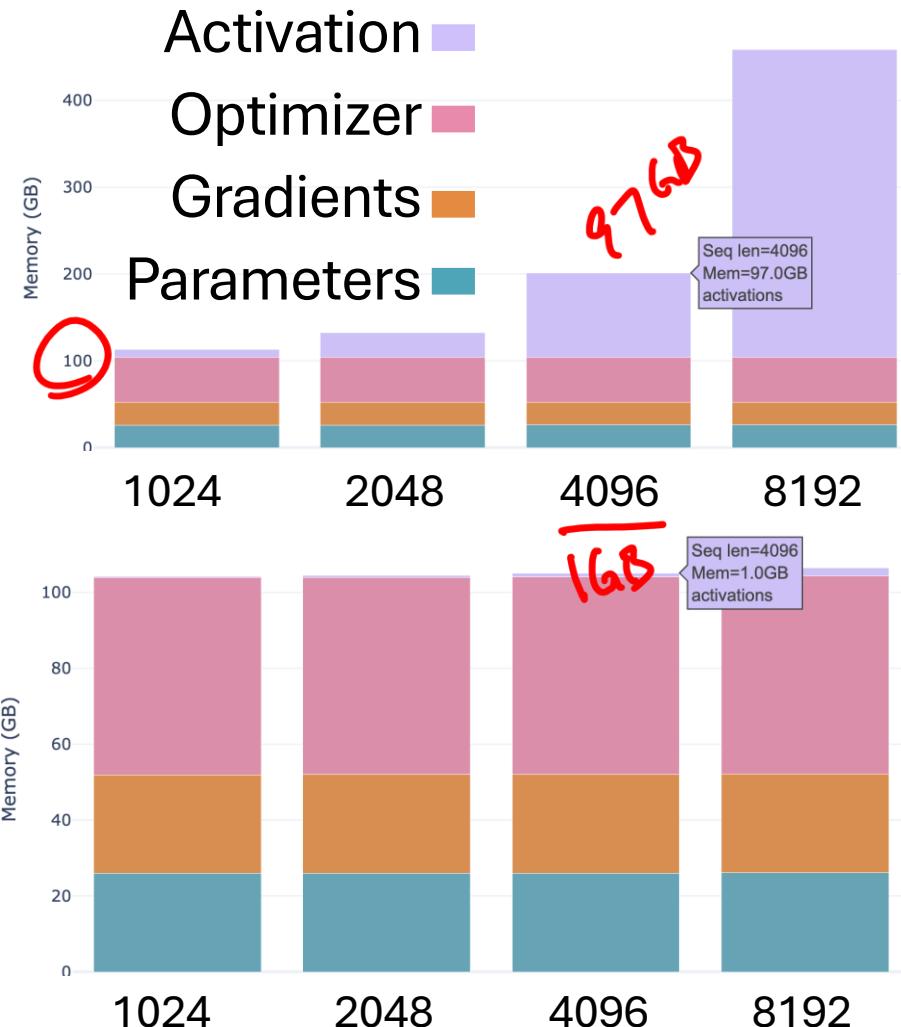
- **Full:** Save at the transition of layers
  - 👍 Saves most memory
  - 👎 Most expensive – 30-40%
- **Selective:** Save MLP; discard attn
  - 70% memory reduction
  - 2.7% compute cost



Korthikanti et al. 2022, Reducing Activation Recomputation in Large Transformer Models



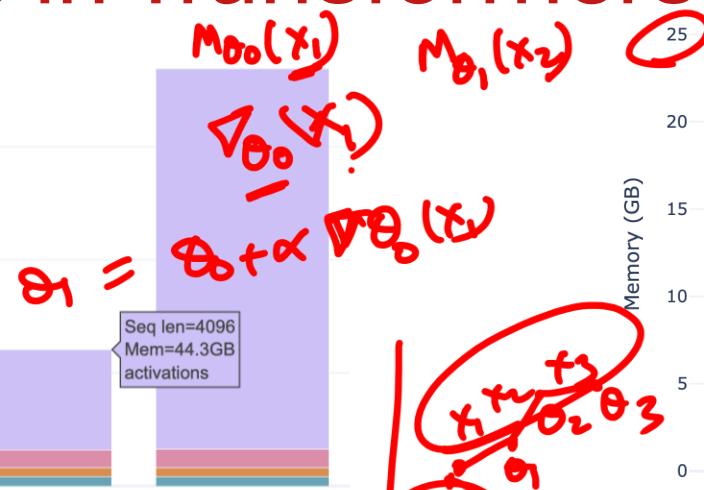
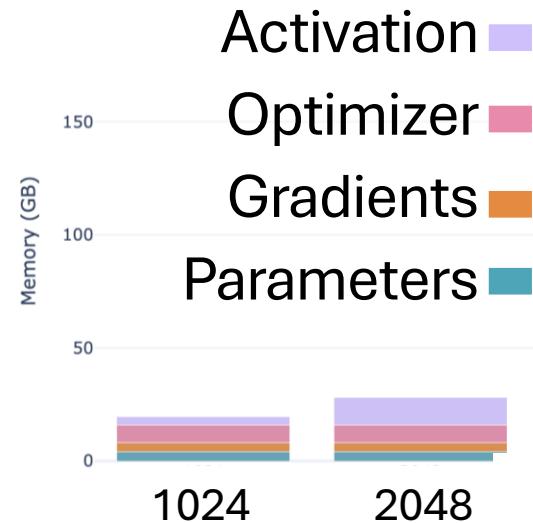
# Memory usage in Transformers – 8B Model , $bs = 1$



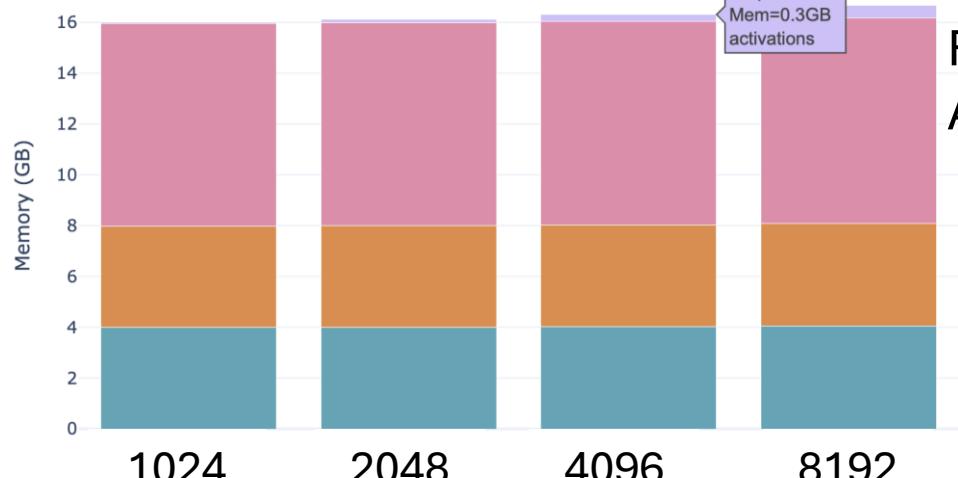
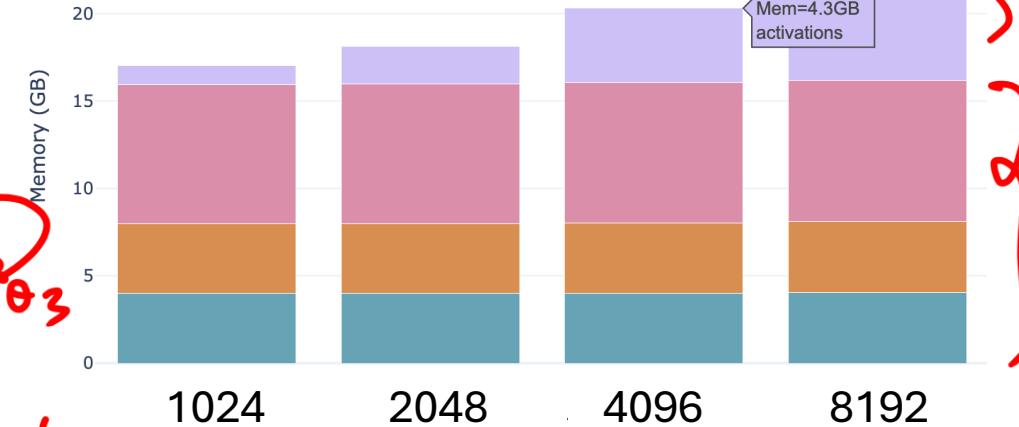
$$\theta_0 \rightarrow \theta_1 \rightarrow \theta_2$$

$$x_1 \quad x_2$$

# Memory usage in Transformers – 1B Model , $bs = 1$



Selective Activation



Full Activation

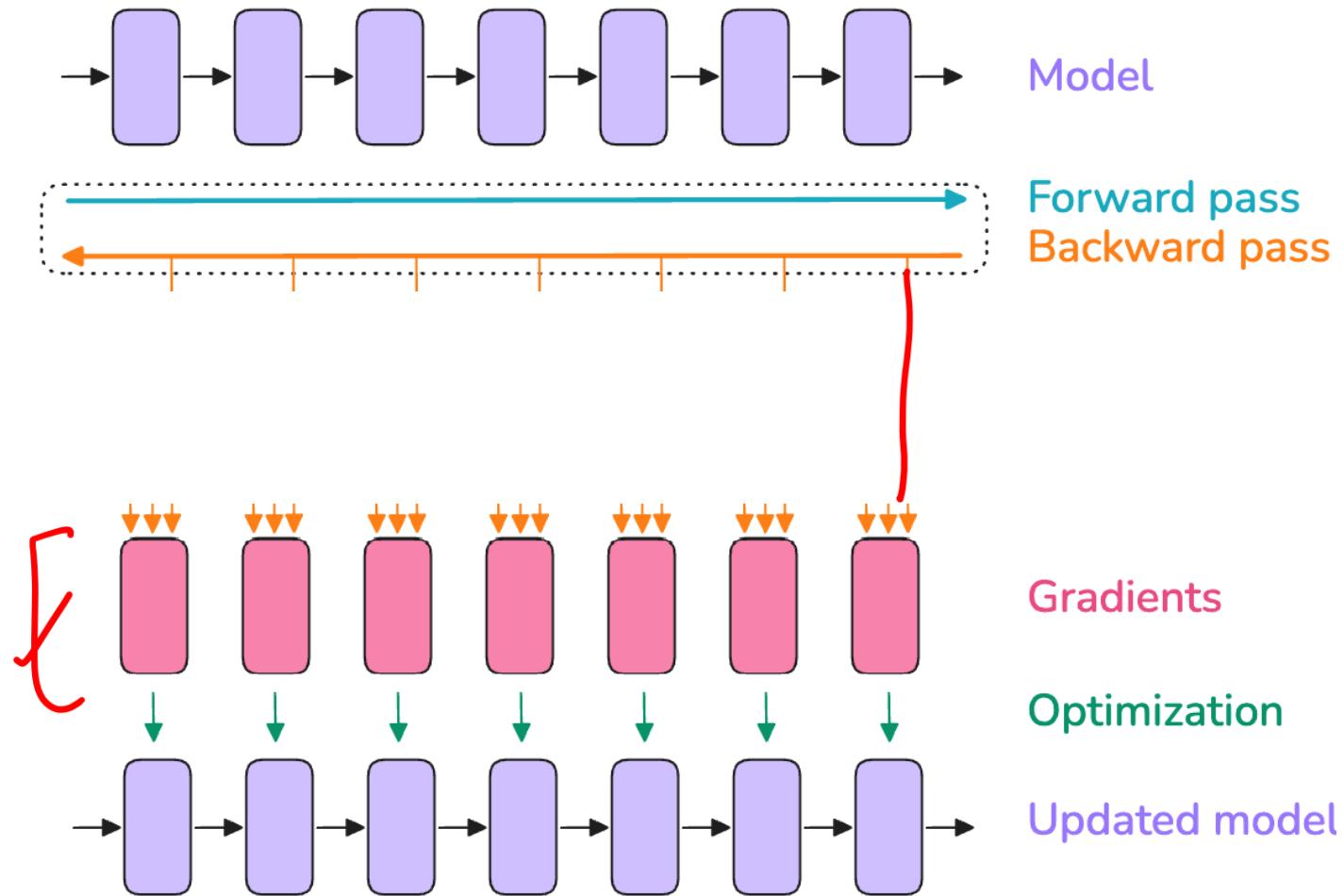
$$m_{act} = L * 4 * \left( \frac{34 * seq * h}{bs} + 5 * n_{leads} * seq^2 * bs \right)$$

Activation memory – grows linearly with batch size

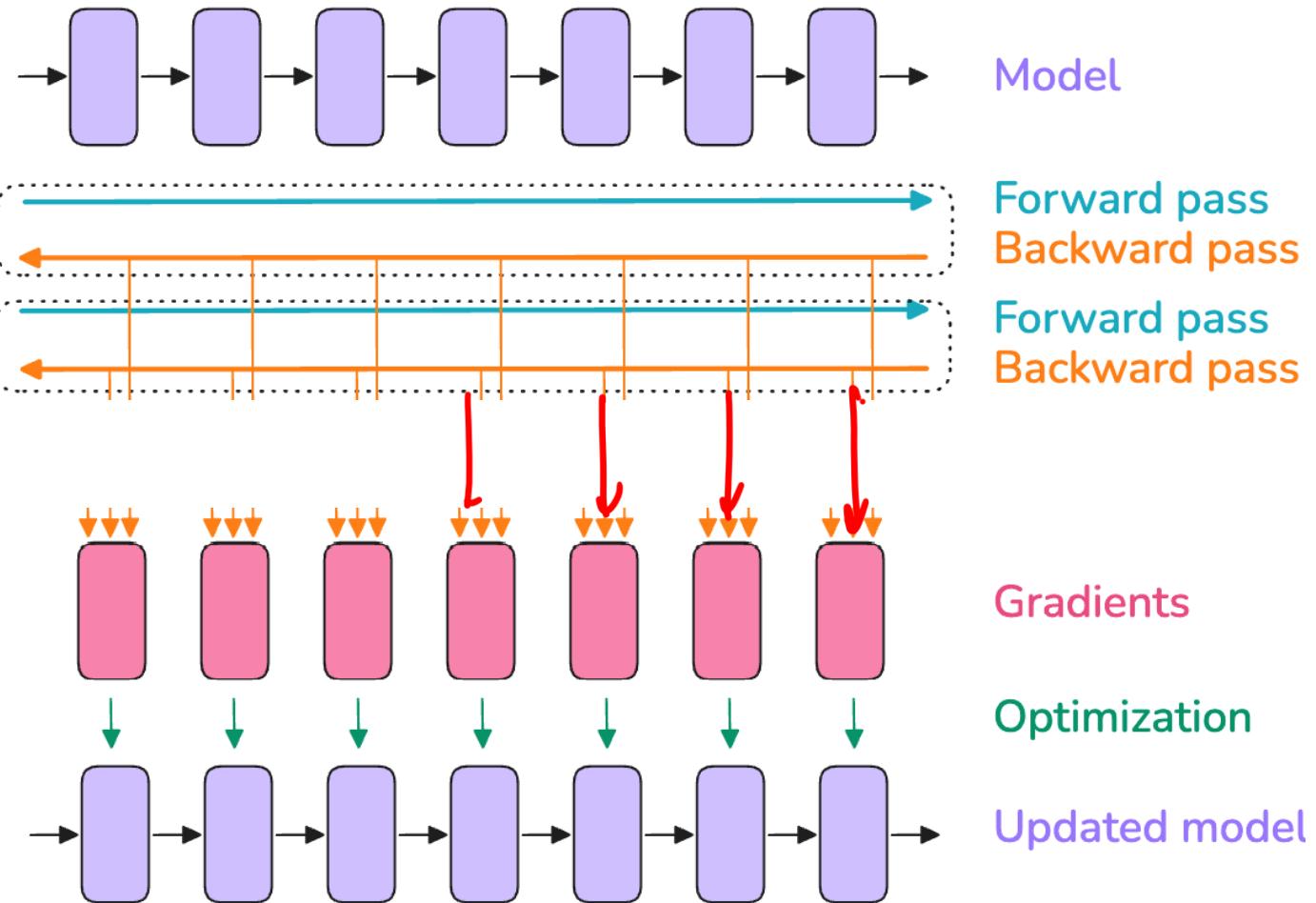
How to increase the batch size?

$$(P_W)_1 + (P_W)_2 + \dots + (P_W)_{25}$$

# How to increase batch size? – *Gradient Accumulation*

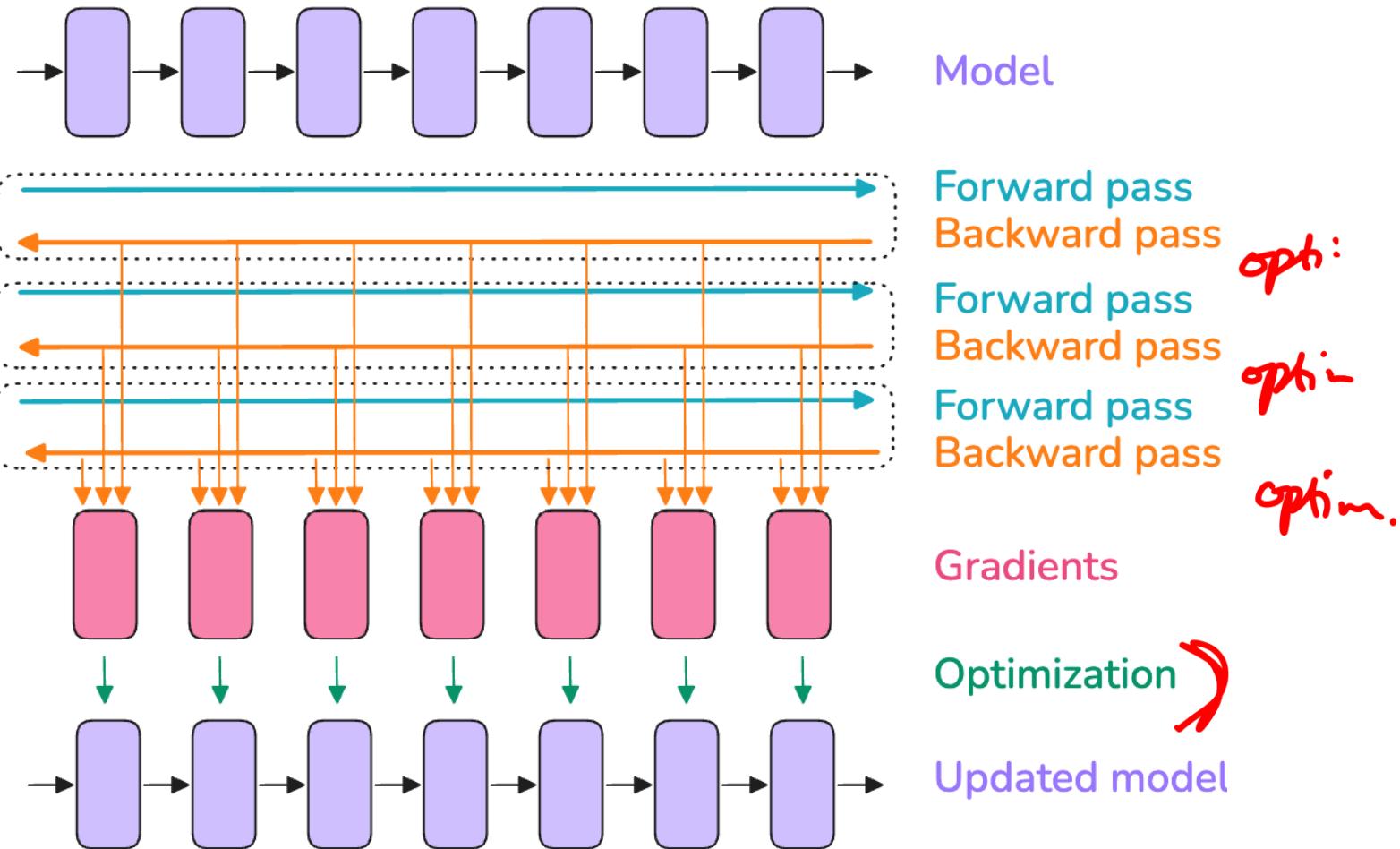


# How to increase batch size? – *Gradient Accumulation*

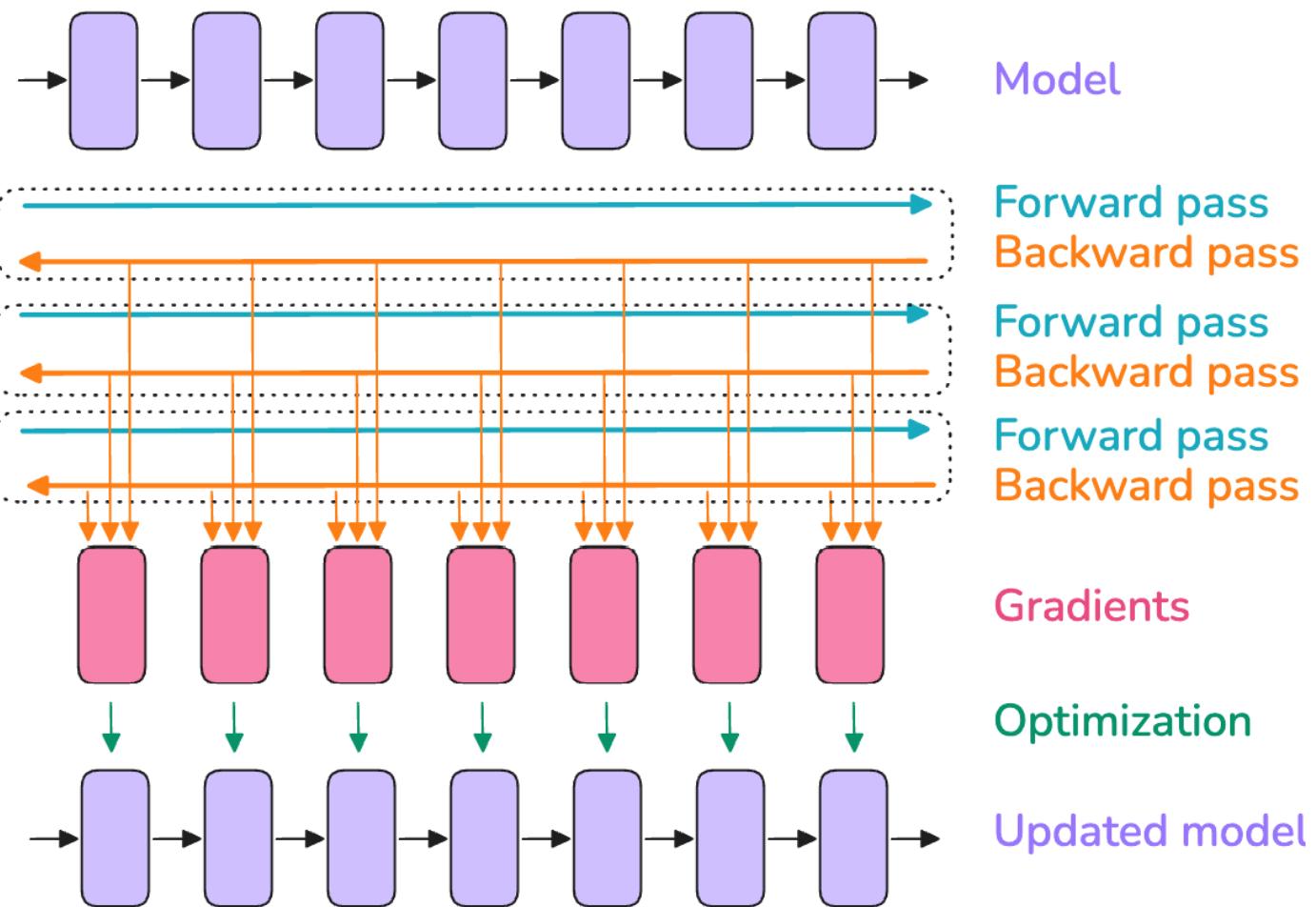


# How to increase batch size? – *Gradient Accumulation*

---



# How to increase batch size? – *Gradient Accumulation*



- Process smaller micro-batches sequentially
- $bs = gbs = mbs * grad\_acc$

Only one micro-batch's worth of activations needs to be kept in memory at a time  
 $(g_1 + g_2 + g_3)$

Requires multiple consecutive forward/backward passes per optimization step

How to speedup?  
softmax with GPU  
 $G_1$        $G_2$