

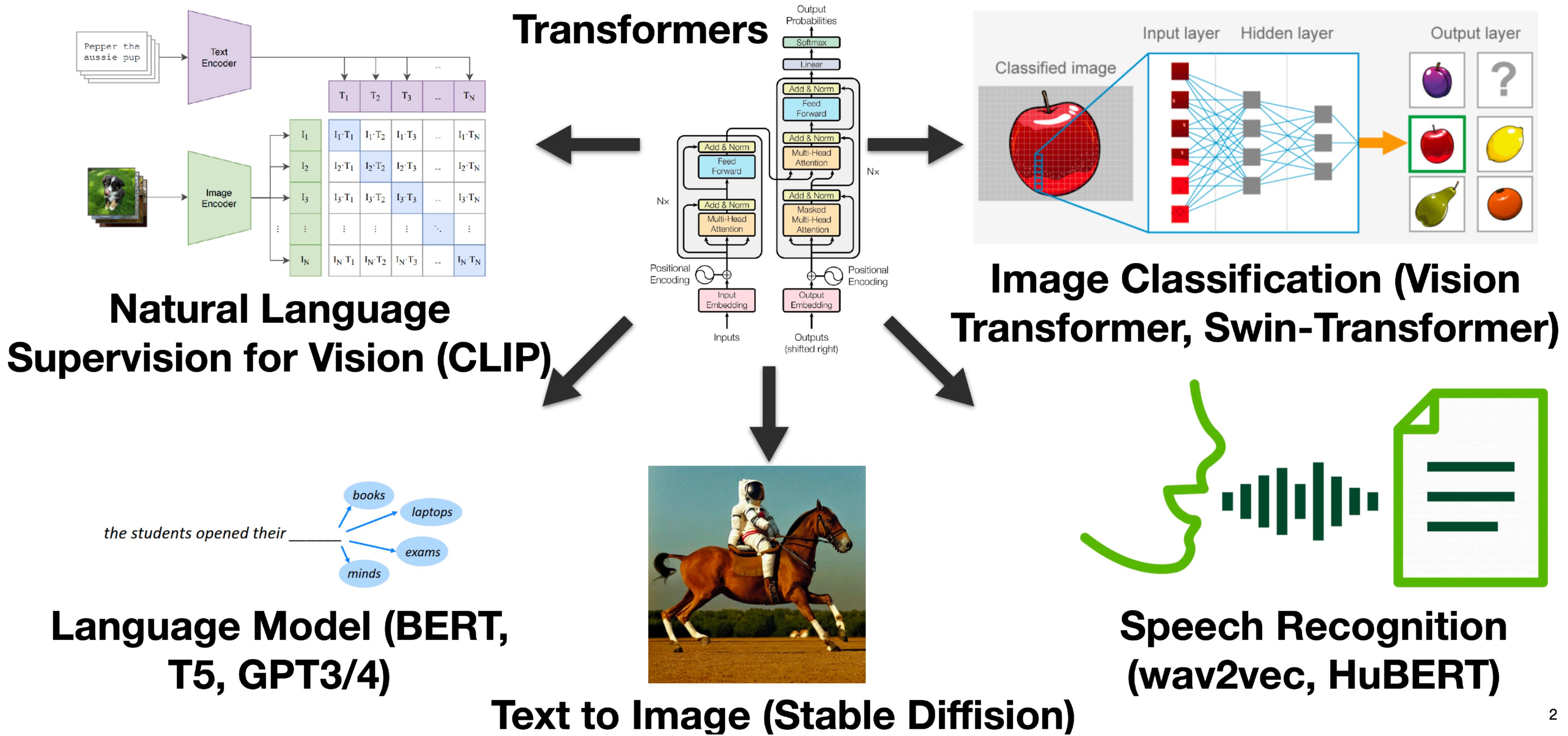


COS 597A

# LI: Overview - Systems for LLM Training and Inference

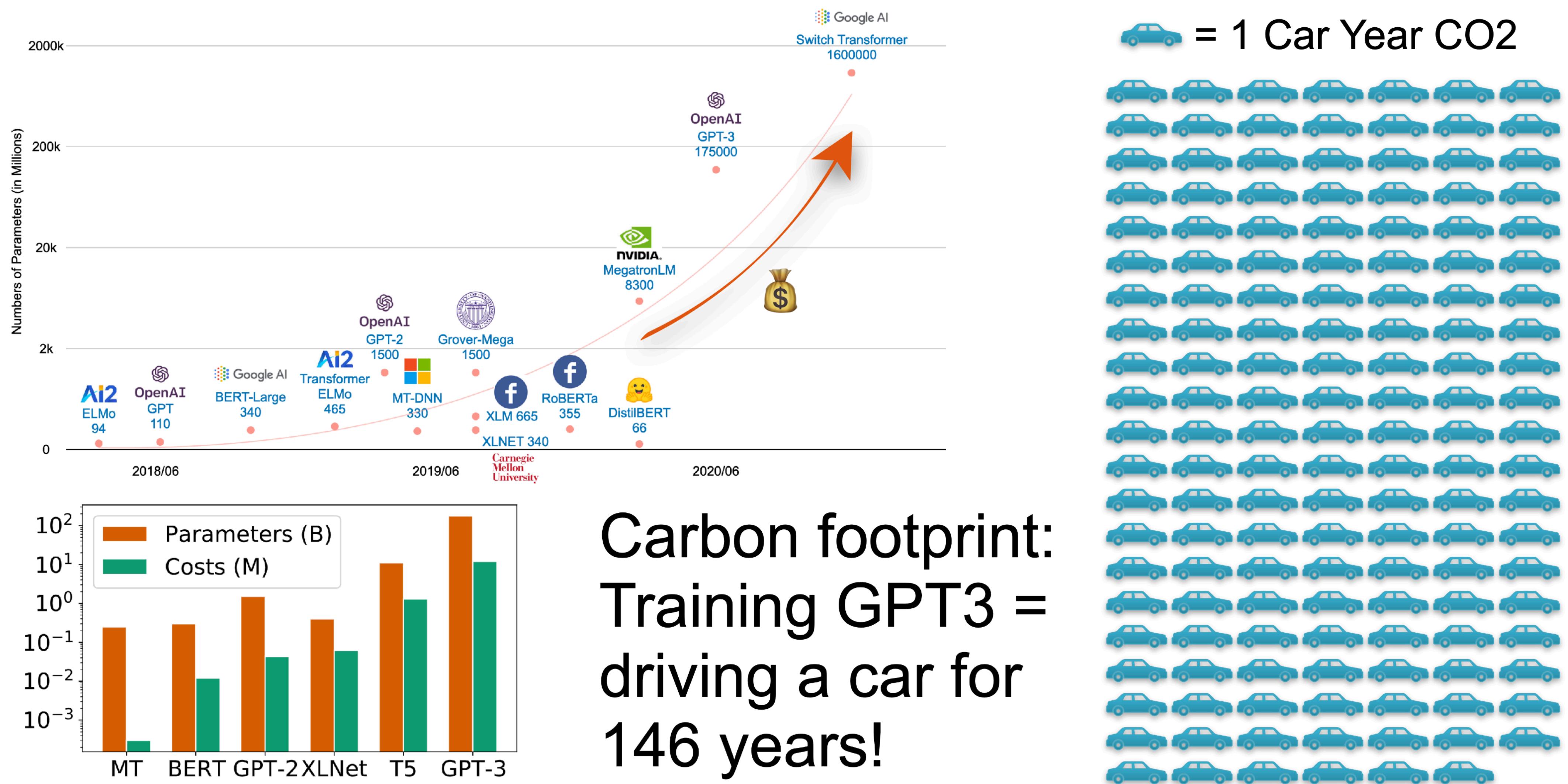
Fall 2025

# Transformer models as universal architecture



Slide credit: Lei Li

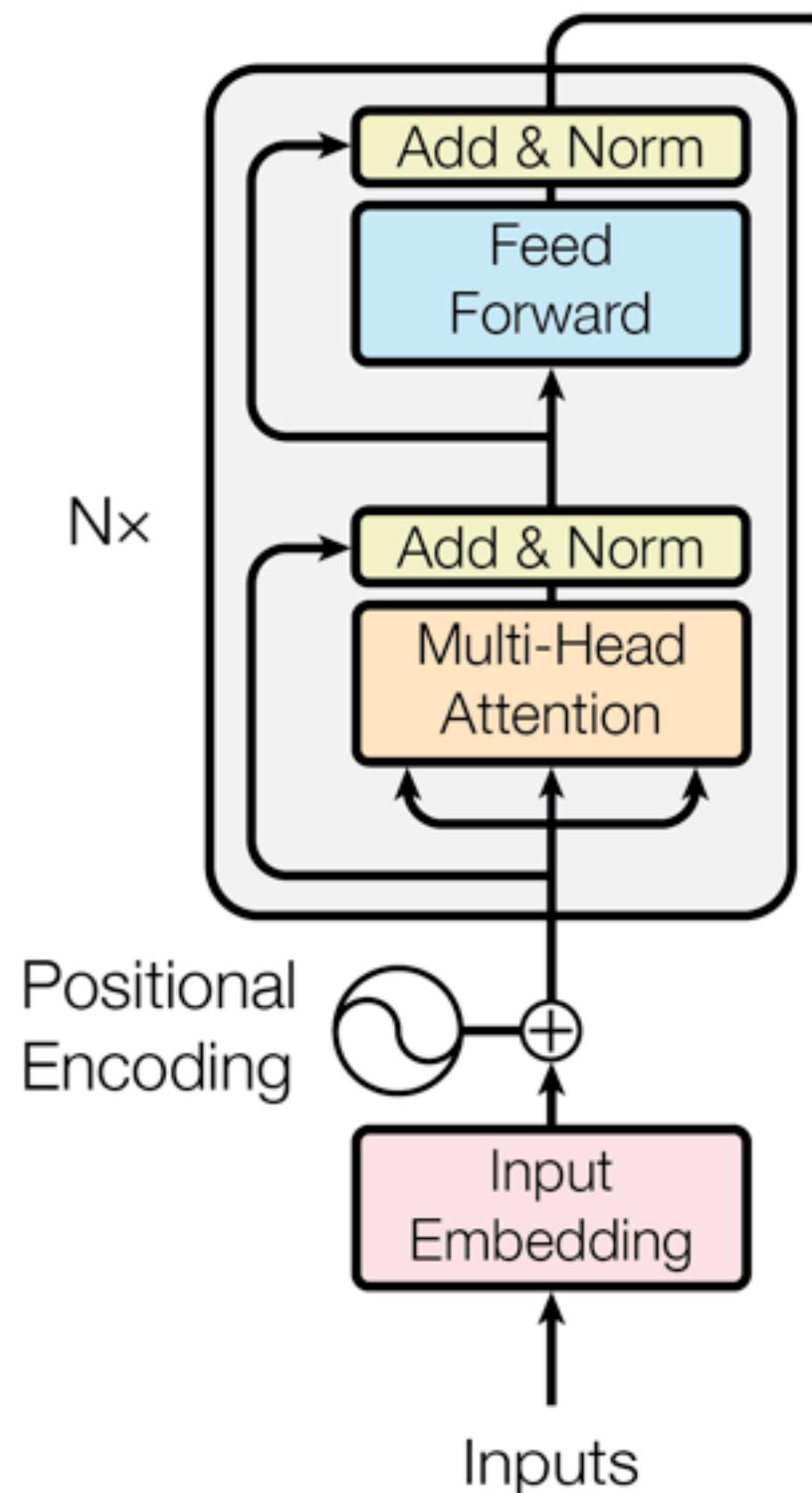
# Why we need systems optimization? LLMs are expensive



# Overview

- (All) Transformer math you need to know
- Training systems
- Inference systems

# Transformer architecture



From the bottom to the top:

- Input embedding
- Positional encoding
- A stack of Transformer encoder layers

Transformer encoder is a stack of  $N$  layers, which consists of two sub-layers:

- Multi-head attention layer
- Feed-forward layer

$$\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^{d_1} \longrightarrow \mathbf{h}_1, \dots, \mathbf{h}_n \in \mathbb{R}^{d_2}$$

# Self-attention

$X \in \mathbb{R}^{n \times d_1}$  (n = input length)

$$Q = XW^Q \quad K = XW^K \quad V = XW^V$$

$$W^Q \in \mathbb{R}^{d_1 \times d_q}, W^K \in \mathbb{R}^{d_1 \times d_k}, W^V \in \mathbb{R}^{d_1 \times d_v}$$

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

$$\begin{aligned} & \text{softmax}\left(\frac{\mathbf{Q} \times \mathbf{K}^T}{\sqrt{d_k}}\right) \mathbf{V} \\ &= \mathbf{H} \end{aligned}$$

# Feed-forward Network (MLP)

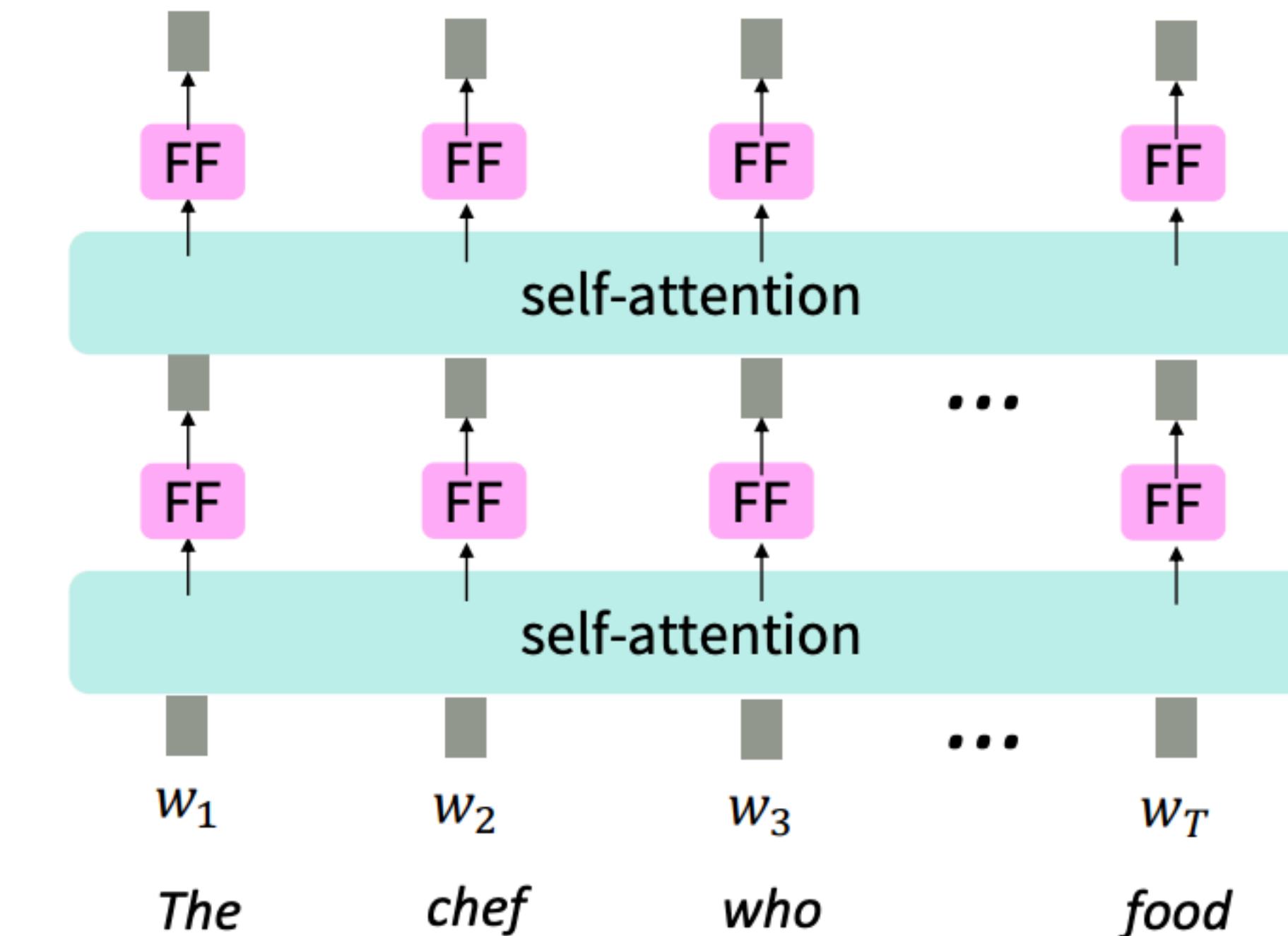
- There are no elementwise nonlinearities in self-attention; stacking more self-attention layers just re-averages value vectors
- Simple fix: add a feed-forward network to post-process each output vector

$$\text{FFN}(\mathbf{x}_i) = \text{ReLU}(\mathbf{x}_i \mathbf{W}_1 + \mathbf{b}_1) \mathbf{W}_2 + \mathbf{b}_2$$

$$\mathbf{W}_1 \in \mathbb{R}^{d \times d_{ff}}, \mathbf{b}_1 \in \mathbb{R}^{d_{ff}}$$

$$\mathbf{W}_2 \in \mathbb{R}^{d_{ff} \times d}, \mathbf{b}_2 \in \mathbb{R}^d$$

In practice, they use  $d_{ff} = 4d$



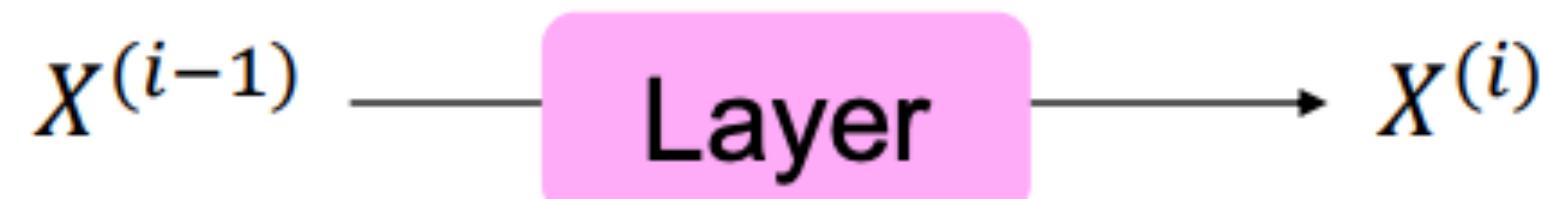
This is actually where the majority of the compute and parameters go!

# Residual connection & layer normalization

Add & Norm:  $\text{LayerNorm}(x + \text{Sublayer}(x))$

## Residual connections (He et al., 2016)

Instead of  $X^{(i)} = \text{Layer}(X^{(i-1)})$  ( $i$  represents the layer)



We let  $X^{(i)} = X^{(i-1)} + \text{Layer}(X^{(i-1)})$ , so we only need to learn “the residual” from the previous layer



Gradient through the residual connection is 1 - good for propagating information through layers

Residual connection unlocks deep networks!

# Residual connection & layer normalization

Add & Norm:  $\text{LayerNorm}(x + \text{Sublayer}(x))$

**Layer normalization** (Ba et al., 2016) helps train model faster

Idea: normalize the hidden vector values to unit mean and stand deviation within each layer

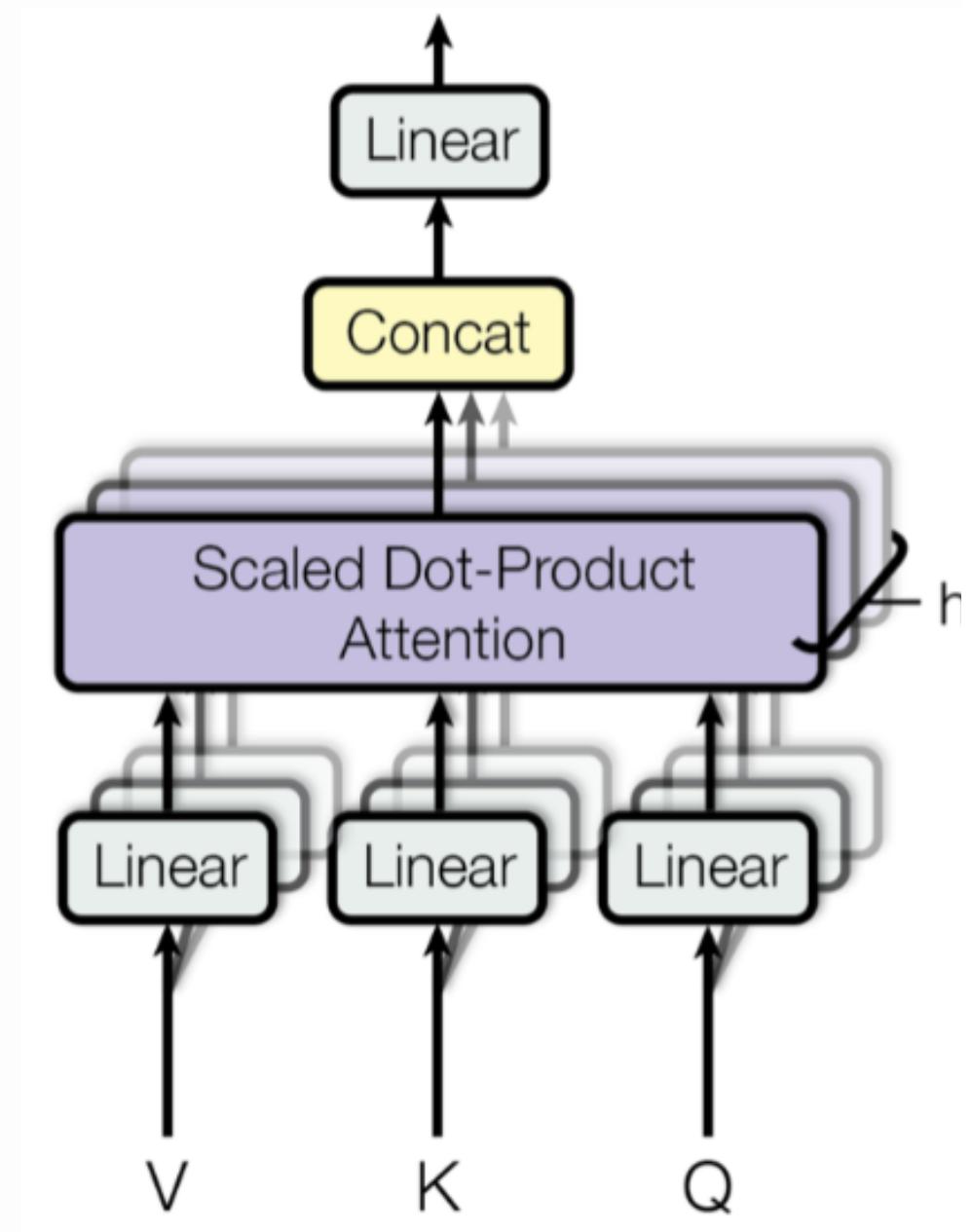
$$y = \frac{x - \mathbb{E}[x]}{\sqrt{\text{Var}[x] + \epsilon}} * \gamma + \beta$$

$\gamma, \beta \in \mathbb{R}^d$  are learnable parameters

LayerNorm is crucial for stable training with deep networks

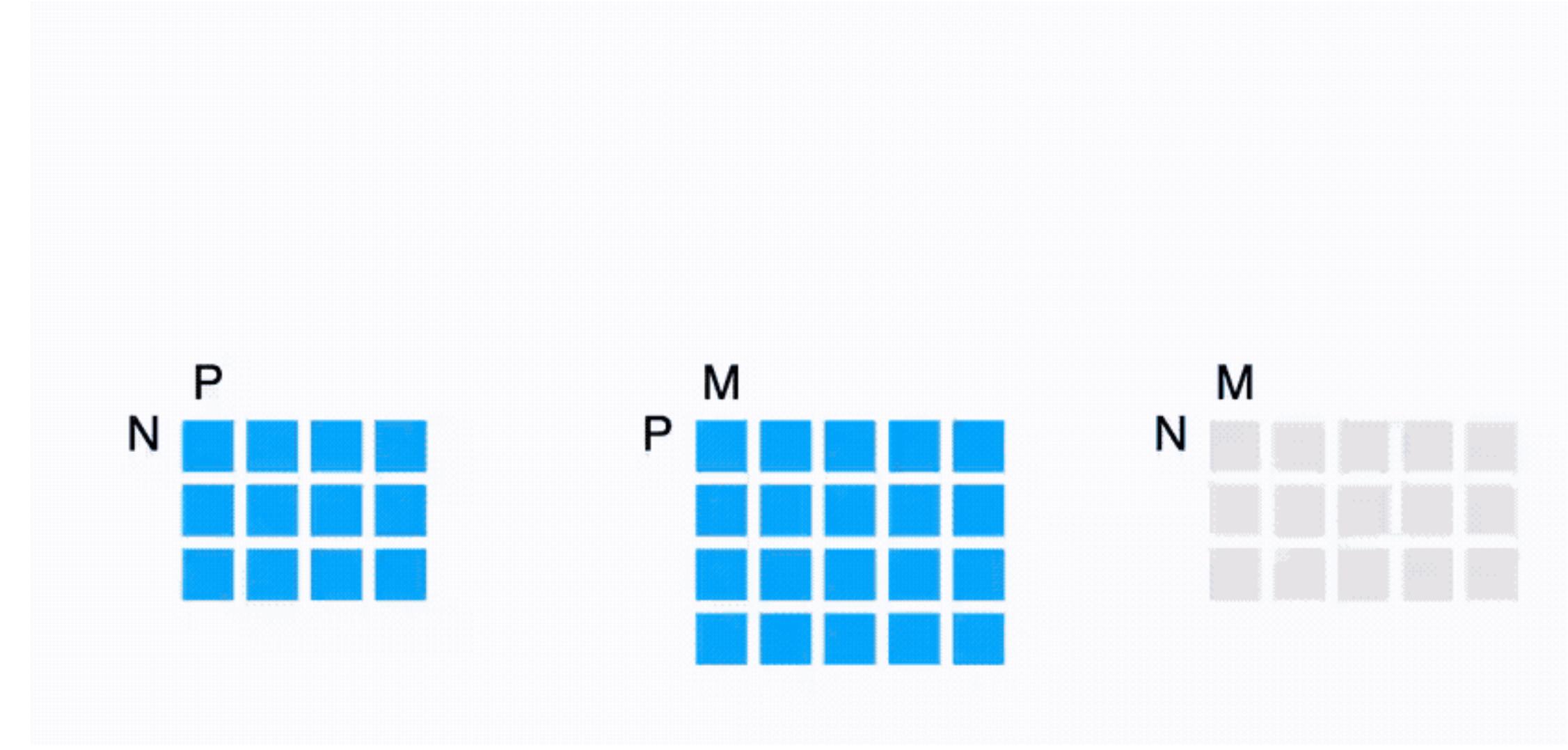
# Transformer Math: Parameters

- Embed dim D, number of layers L, vocab size V, context length N
- MHA layer:  $4D^2$  ( $3D^2$  for QKV projection,  $D^2$  for output projection)
- FFN layer:  $8D^2$   
 $\text{FFN}(\mathbf{x}_i) = W_2\phi(W_1\mathbf{x}_i + \mathbf{b}_1) + \mathbf{b}_2$   
 $\mathbf{W}_1 \in \mathbb{R}^{d \times d_{ff}}, \mathbf{W}_2 \in \mathbb{R}^{d_{ff} \times d}$   
Typically  $d_{ff} = 4d$
- Embedding and LM head:  $DV$  each
- Total:  $12LD^2 + DV$  if shared embedding,  $12LD^2 + 2DV$  if not shared



# Transformer Math: Matrix Multiply

- Matmul size  $(N \times P) @ (P \times M)$  has  $2MNP$  FLOPS



- Per input vector, each weight matrix size  $N \times P$  has  $2NP$  FLOPS
- Forward pass FLOPS per input vector:  $2 \times$  no. (non-embedding) params

# Transformer Math: Forward and Backward FLOPS

- Forward:  $A = W X$
- Backward:
  - Weight gradient  $dW = dA X^T$
  - Activation gradient  $dX = W^T dA$
  - Twice the FLOPS of forward
- Forward + backward FLOPS per input vector:  $6 \times \text{no. (non-embedding) params}$

# Transformer Math: Attention FLOPS

- Forward, for input length N:  $4N^2D$  ( $2N^2D$  for  $QK^T$ ,  $2N^2D$  for  $PV$ )  
N: context length, D: model dimension  
Per input vector:  $4ND$

- Backward: 2x forward
- Forward + backward FLOPS per input vector:  $12LND$   
L: number of layers

	S1	S1	EOT	S2	S2
S1	1	0	0	0	0
S1	1	1	0	0	0
EOT	1	1	1	0	0
S2	1	1	1	1	0
S2	1	1	1	1	1

- For causal attention: only compute half the entries ->  $6LND$   
Convention varies on how to count  
(should it be  $12LND$ ,  $6LND$ , or even  $8LND$  or  $7LND$  due to recompute in the backward pass?)

# Transformer Math: Total FLOPS

- Forward + backward FLOPS per input vector:  
 $6 \times \text{no. (non-embedding) params} + 12LND$   
L layers, context length N, model dim D
- Typically approximated as:  $6 \times \text{no. (non-embedding) params}$  when context length is not too long
- Total FLOPS when trained on T tokens:  $\approx 6 \times \text{no. (non-embedding) params} \times \text{no. tokens}$

# Hardware

## NVIDIA A100 TENSOR CORE GPU SPECIFICATIONS (SXM4 AND PCIE FORM FACTORS)

	A100 80GB PCIe	A100 80GB SXM
FP64	<b>9.7 TFLOPS</b>	
FP64 Tensor Core		<b>19.5 TFLOPS</b>
FP32	<b>19.5 TFLOPS</b>	
Tensor Float 32 (TF32)	<b>156 TFLOPS   312 TFLOPS*</b>	
BFLOAT16 Tensor Core	<b>312 TFLOPS   624 TFLOPS*</b>	
FP16 Tensor Core	<b>312 TFLOPS   624 TFLOPS*</b>	
INT8 Tensor Core	<b>624 TOPS   1248 TOPS*</b>	
GPU Memory	<b>80GB HBM2e</b>	<b>80GB HBM2e</b>
GPU Memory Bandwidth	<b>1,935GB/s</b>	<b>2,039GB/s</b>
Max Thermal Design Power (TDP)	<b>300W</b>	<b>400W***</b>
Multi-Instance GPU	<b>Up to 7 MIGs @ 10GB</b>	<b>Up to 7 MIGs @ 10GB</b>
Form Factor	<b>PCIe dual-slot air cooled or single-slot liquid cooled</b>	<b>SXM</b>
Interconnect	<b>NVIDIA® NVLink® Bridge for 2 GPUs: 600GB/s **</b> <b>PCIe Gen4: 64GB/s</b>	<b>NVLink: 600GB/s</b> <b>PCIe Gen4: 64GB/s</b>
Server Options	<b>Partner and NVIDIA-Certified Systems™ with 1-8 GPUs</b>	<b>NVIDIA HGX™ A100- Partner and NVIDIA-Certified Systems with 4, 8, or 16 GPUs</b> <b>NVIDIA DGX™ A100 with 8 GPUs</b>

\* With sparsity

## Technical Specifications

<b>H100 SXM</b>	
<b>FP64</b>	34 teraFLOPS
<b>FP64 Tensor Core</b>	67 teraFLOPS
<b>FP32</b>	67 teraFLOPS
<b>TF32 Tensor Core*</b>	989 teraFLOPS
<b>BFLOAT16 Tensor Core*</b>	1,979 teraFLOPS
<b>FP16 Tensor Core*</b>	1,979 teraFLOPS
<b>FP8 Tensor Core*</b>	3,958 teraFLOPS
<b>INT8 Tensor Core*</b>	3,958 TOPS
<b>GPU Memory</b>	80GB
<b>GPU Memory Bandwidth</b>	3.35TB/s
<b>Decoders</b>	7 NVDEC 7 JPEG
<b>Max Thermal Design Power (TDP)</b>	Up to 700W (configurable)
<b>Multi-Instance GPUs</b>	Up to 7 MIGs @ 10GB each
<b>Form Factor</b>	SXM
<b>Interconnect</b>	NVIDIA NVLink™: 900GB/s PCIe Gen5: 128GB/s
<b>Server Options</b>	NVIDIA HGX H100 Partner and NVIDIA-Certified Systems™ with 4 or 8 GPUs NVIDIA DGX H100 with 8 GPUs
<b>NVIDIA Enterprise</b>	Add-on

\*With sparsity

**BF16 Dense:  
989 TFLOPS**

# How long does it take to train a model?

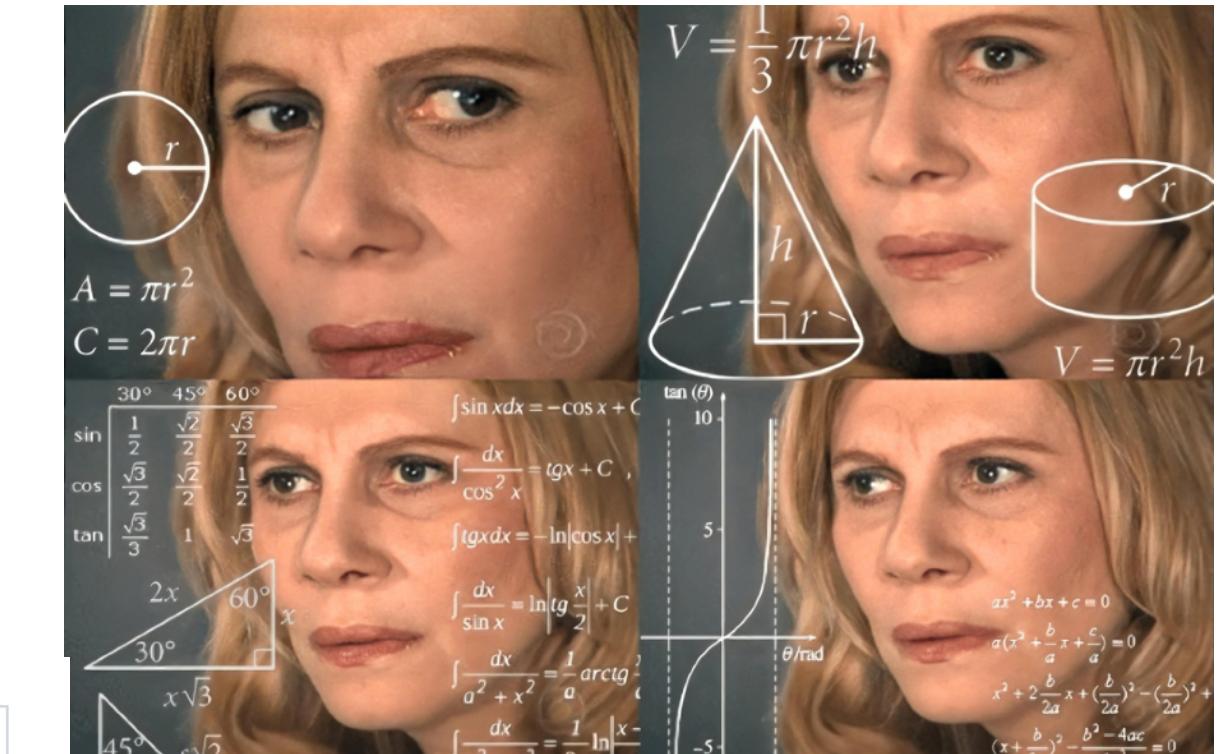
- Model FLOPS utilization (MFU)  
MFU = Model FLOPS per sec / theoretical max TFLOPS per sec
- Typical MFU: 30-50% (e.g. 300-500 TFLOPS/sec per H100)
- Why: memory-bandwidth bound operations, communication, power-throttling
- E.g., how many H100 hours does it take to train a X billion model to Y trillion tokens?  
 $\approx 6 * 10^9 X * 10^{12} Y \text{ FLOPS} / (400 * 10^{12} * 3600 \text{ sec})$

# How long does it take to train a Llama?

- How many H100 hours does it take to train a X billion model to Y trillion tokens?  
 $\approx 6 * 10^9 X * 10^{12} Y \text{ FLOPS} / (400 * 10^{12} * 3600 \text{ sec})$

- E.g., Llama-3: 405B, 15T tokens  
 $\approx 6 * 405e9 * 15e12 / (400e12 * 3600)$   
 $= 25.3\text{M hours} = 66 \text{ days on 16K H100}$

	Training Time (GPU hours)	Training Power Consumption (W)	Training Location-Based Greenhouse Gas Emissions (tons CO2eq)	Training Market-Based Greenhouse Gas Emissions (tons CO2eq)
Llama 3.1 8B	1.46M	700	420	0
Llama 3.1 70B	7.0M	700	2,040	0
Llama 3.1 405B	30.84M	700	8,930	0
Total	39.3M		11,390	0



# Overview

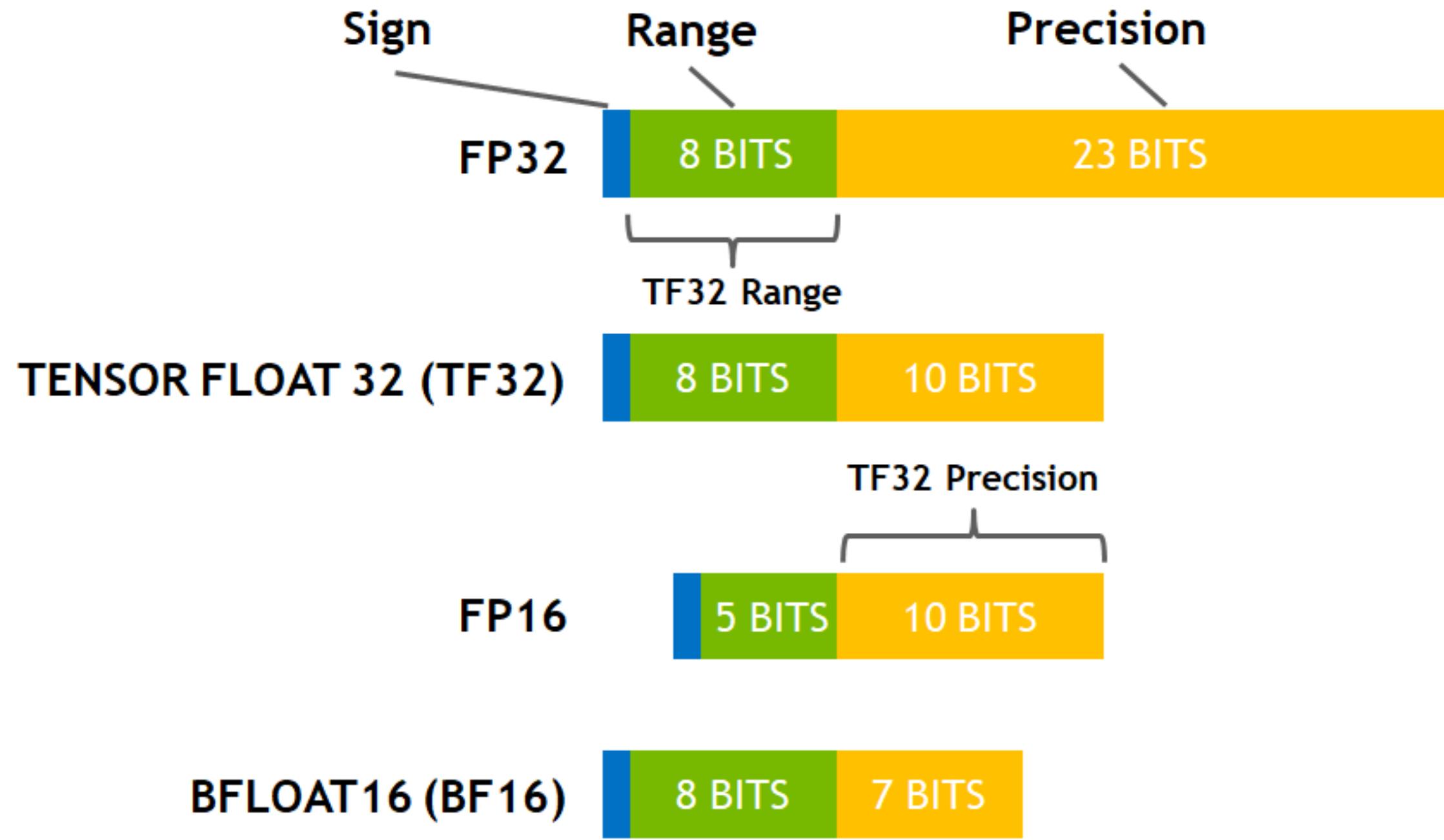
- (All) Transformer math you need to know
- **Training systems**
- Inference systems

# Training: Mixed Precision

Technical Specifications	
<b>H100 SXM</b>	
<b>FP64</b>	34 teraFLOPS
<b>FP64 Tensor Core</b>	67 teraFLOPS
<b>FP32</b>	67 teraFLOPS
<b>TF32 Tensor Core*</b>	989 teraFLOPS
<b>BFLOAT16 Tensor Core*</b>	1,979 teraFLOPS
<b>FP16 Tensor Core*</b>	1,979 teraFLOPS
<b>FP8 Tensor Core*</b>	3,958 teraFLOPS
<b>INT8 Tensor Core*</b>	3,958 TOPS
<b>GPU Memory</b>	80GB
<b>GPU Memory Bandwidth</b>	3.35TB/s
<b>Decoders</b>	7 NVDEC 7 JPEG
<b>Max Thermal Design Power (TDP)</b>	Up to 700W (configurable)
<b>Multi-Instance GPUs</b>	Up to 7 MIGs @ 10GB each
<b>Form Factor</b>	SXM
<b>Interconnect</b>	NVIDIA NVLink™: 900GB/s  PCIe Gen5: 128GB/s
<b>Server Options</b>	NVIDIA HGX H100 Partner and NVIDIA-Certified Systems™ with 4 or 8 GPUs  NVIDIA DGX H100 with 8 GPUs
<b>NVIDIA Enterprise</b>	Add-on

\*With sparsity

**BF16 Dense:  
989 TFLOPS**



**BF16/FP16 FLOPS are much higher than FP32!**

# Training: Automatic Mixed Precision

FP16

- GEMMs + Convolutions can use Tensor Cores
- Most pointwise ops (e.g. add, multiply):  
1/2X memory storage for intermediates,  
2X memory throughput

FP32

- Weight updates benefit from precision
- Loss functions (often reductions) benefit from precision and range
- Softmax, norms, some other ops benefit from precision and range

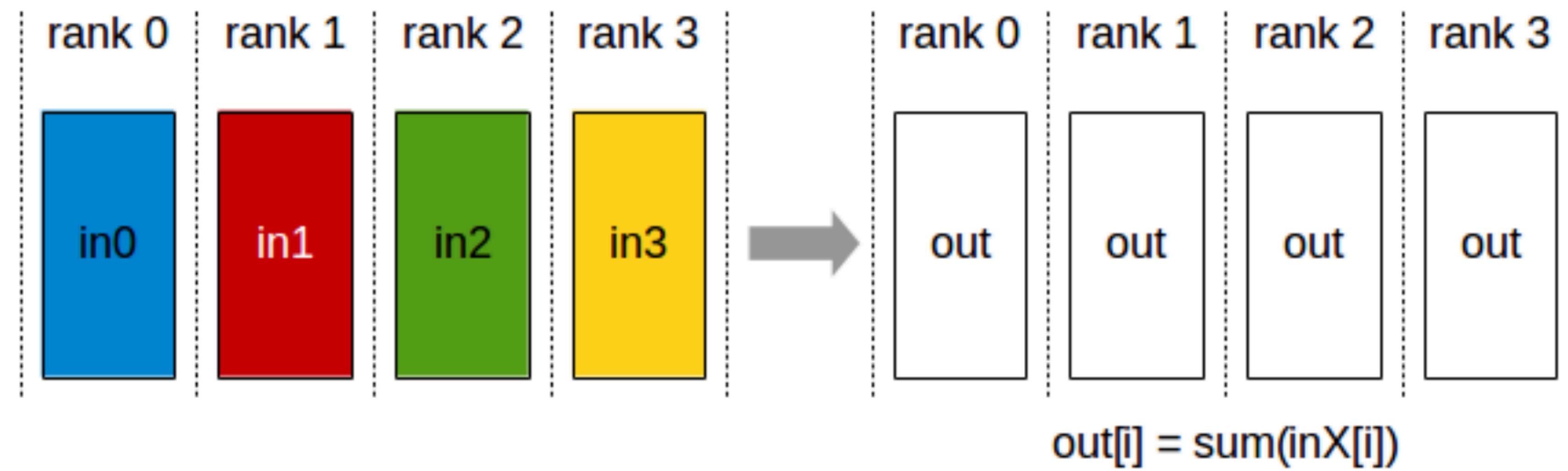


Done automatically if you use PyTorch AMP

# Distributed Training: Communication

## AllReduce

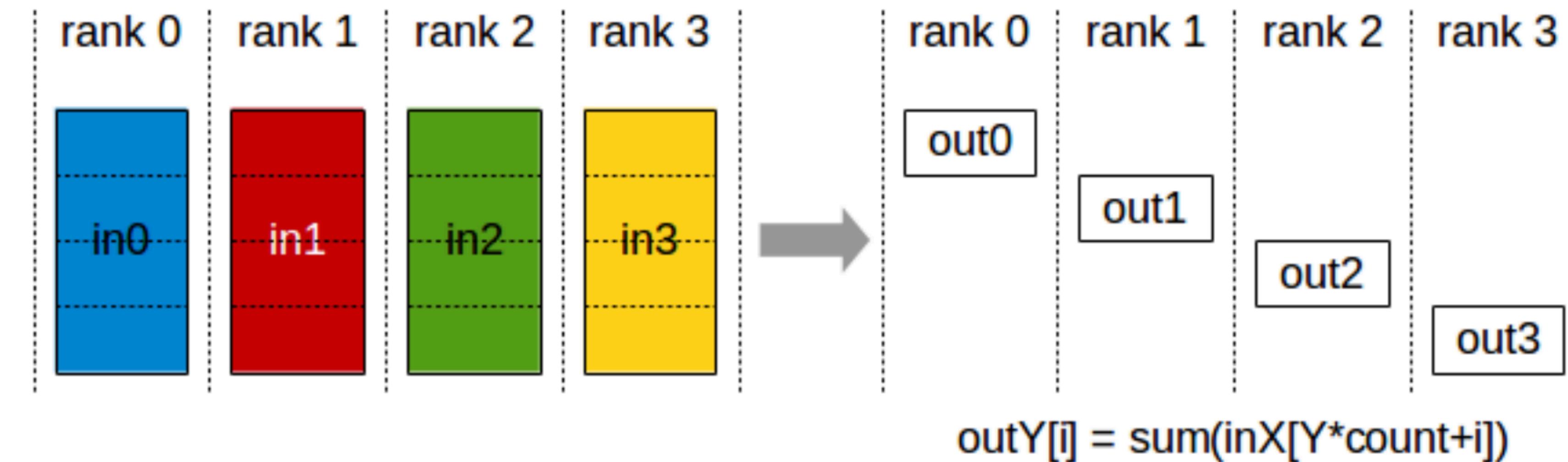
- Compute reduction (sum, min, max) across devices and writing the result in the receive buffers of every rank.



# Distributed Training: Communication

## ReduceScatter

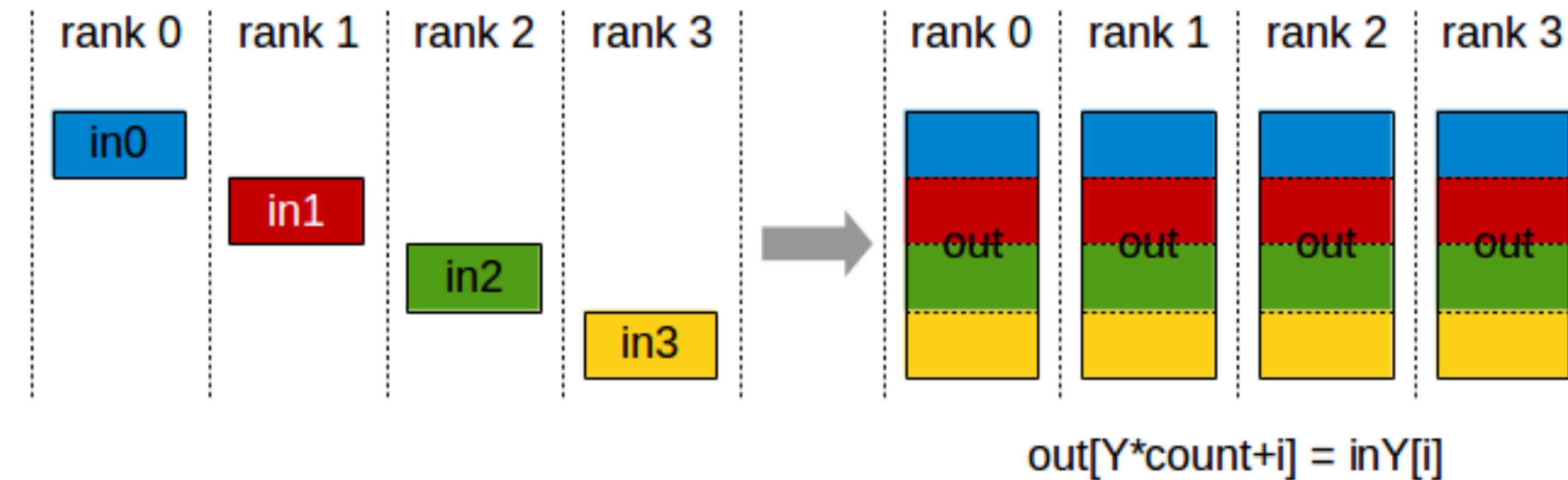
- Compute reduction (sum, min, max) and writing parts of results scattered in ranks



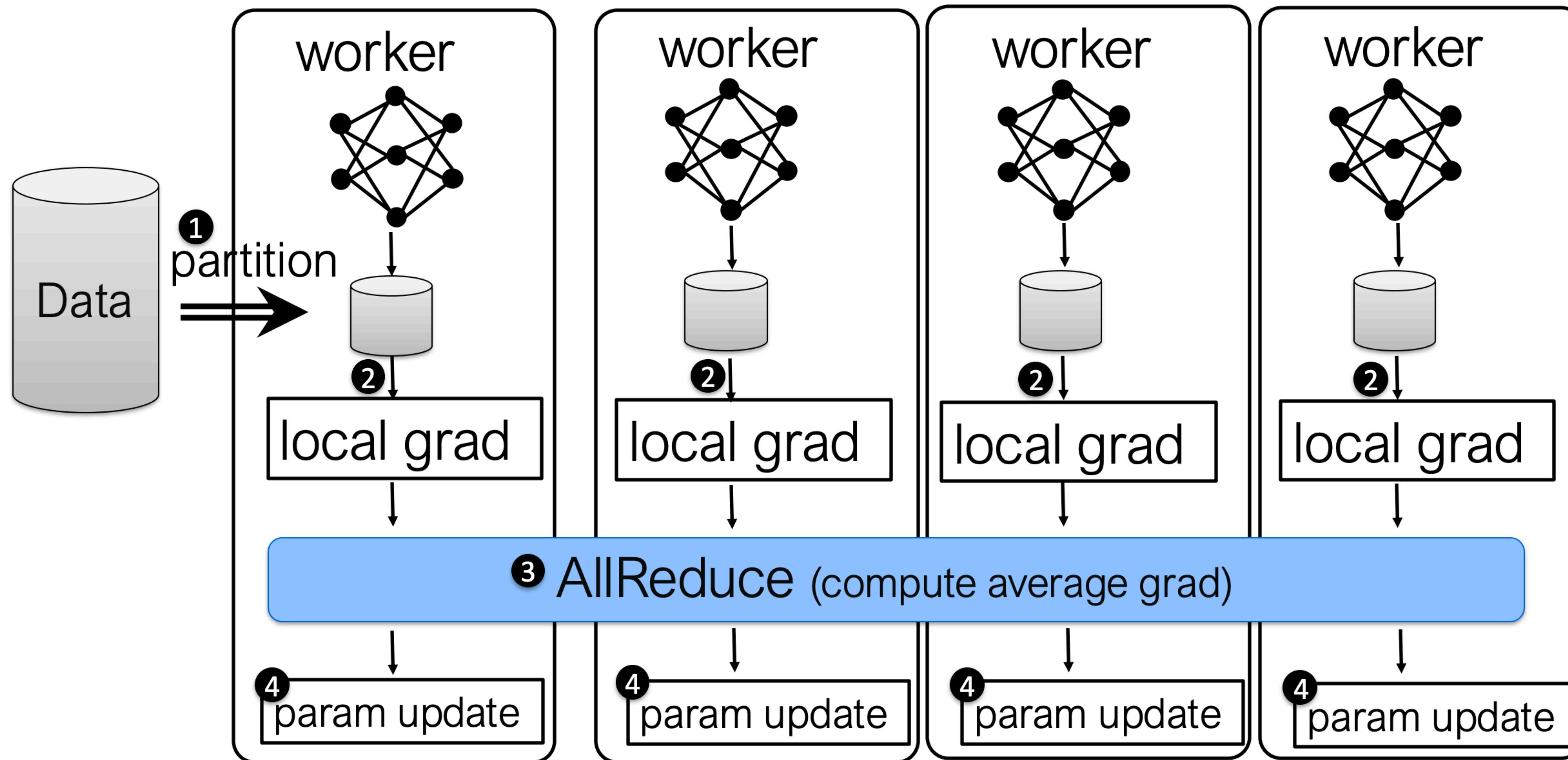
# Distributed Training: Communication

## AllGather

- gathers N values from k ranks into an output of size  $k*N$ , and distributes that result to all ranks (devices).

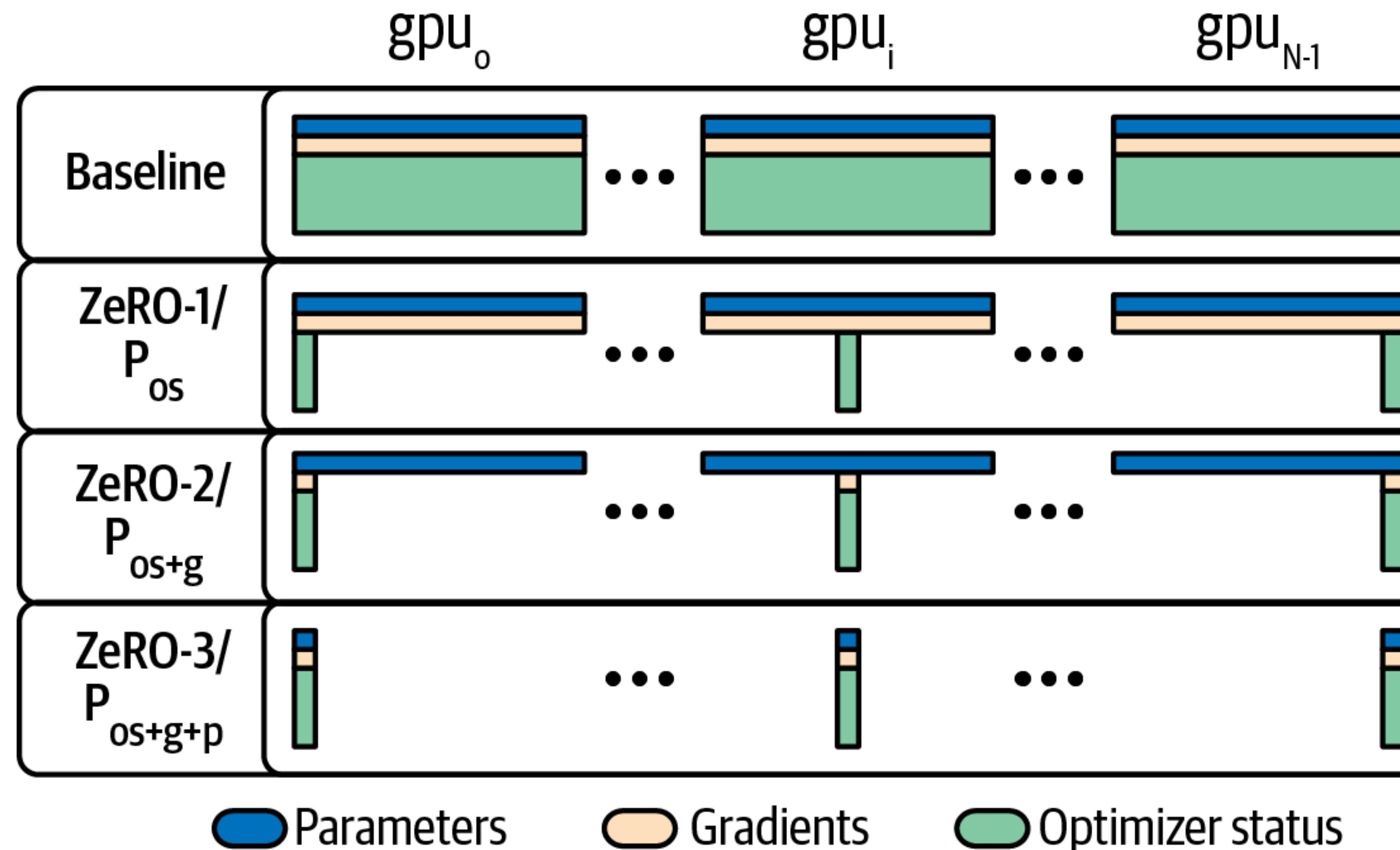


# Distributed Training: Data Parallel



PyTorch calls this Distributed Data Parallel (DDP)

# Distributed Training: Zero Redundancy Optimizer (ZeRO)



PyTorch: Fully Sharded Data Parallel (FSDP)

# Distributed Training: Tensor Parallel

$$\begin{matrix} \text{C} \\ \hline \end{matrix} = \begin{matrix} \text{A} \\ \hline \end{matrix} \times \begin{matrix} \text{B} \\ \hline \end{matrix}$$

Non-distributed

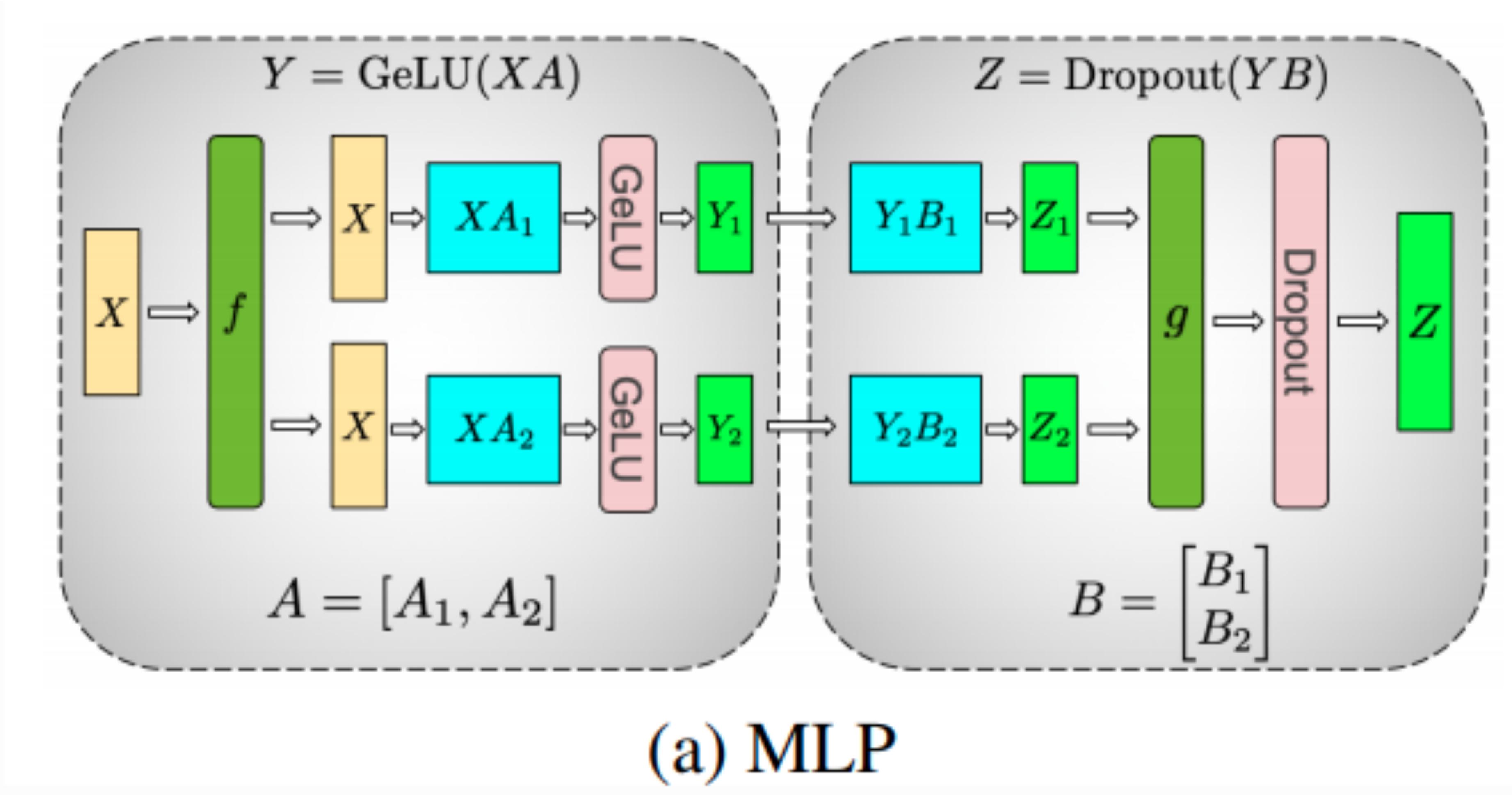
$$\begin{matrix} \text{C} \\ \hline \end{matrix} = \begin{matrix} \text{A} \\ \hline \end{matrix} \times \begin{matrix} \text{B} \\ \hline \end{matrix}$$

all-gather  
along column

$$\begin{matrix} \text{C} \\ \hline \end{matrix} = \begin{matrix} \text{A} \\ \hline \end{matrix} \times \begin{matrix} \text{B} \\ \hline \end{matrix}$$

Column-Splitting Tensor Parallel

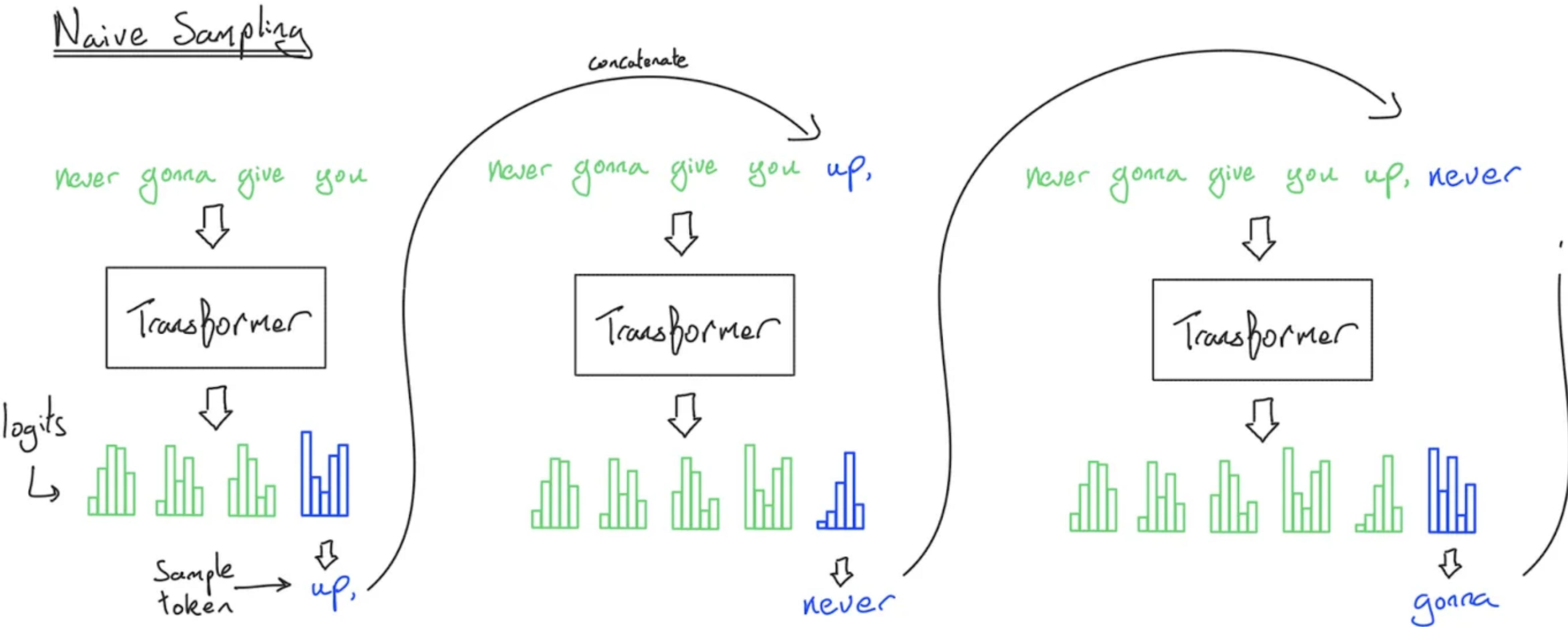
# Distributed Training: Tensor Parallel



# Overview

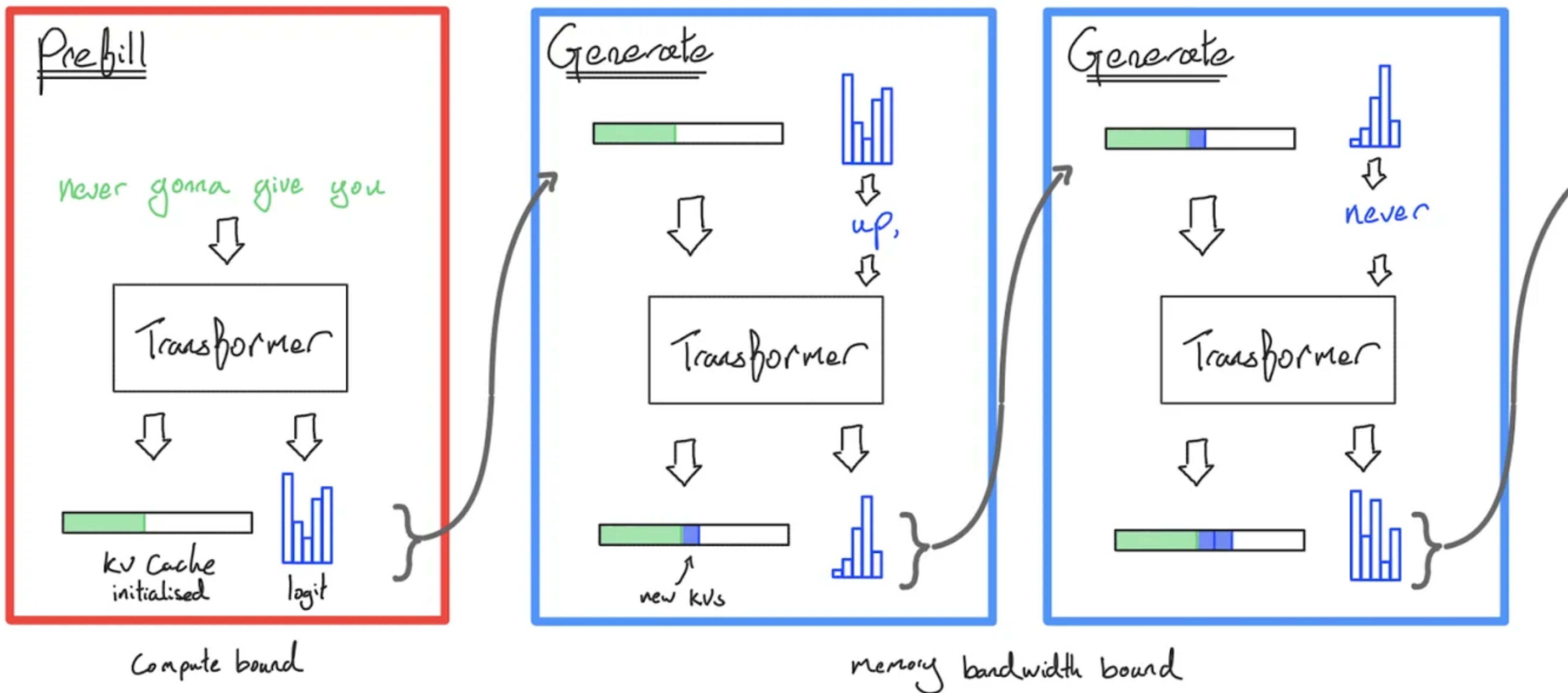
- (All) Transformer math you need to know
- Training systems
- **Inference systems**

# Inference: KV caching



# Inference: KV caching

Sampling with KV cache



# Inference: KV caching

$$\begin{array}{c}
 \left[ \begin{array}{c} Q \\ \vdots \\ \text{q1} \\ \text{q2} \\ \text{q3} \\ \text{q4} \end{array} \right] \times \begin{array}{c} K^T \\ \vdots \\ \text{k1} | \text{k2} | \text{k3} | \text{k4} \end{array} = \begin{array}{c} QK^T \\ \vdots \\ \text{q1}\cdot\text{k1} | \text{q1}\cdot\text{k2} | \text{q1}\cdot\text{k3} | \text{q1}\cdot\text{k4} \\ \text{q2}\cdot\text{k1} | \text{q2}\cdot\text{k2} | \text{q2}\cdot\text{k3} | \text{q2}\cdot\text{k4} \\ \text{q3}\cdot\text{k1} | \text{q3}\cdot\text{k2} | \text{q3}\cdot\text{k3} | \text{q3}\cdot\text{k4} \\ \text{q4}\cdot\text{k1} | \text{q4}\cdot\text{k2} | \text{q4}\cdot\text{k3} | \text{q4}\cdot\text{k4} \end{array} \right] \times \begin{array}{c} V \\ \vdots \\ \text{v1} \\ \text{v2} \\ \text{v3} \\ \text{v4} \end{array} = \begin{array}{c} A \\ \vdots \\ \text{a1} \\ \text{a2} \\ \text{a3} \\ \text{a4} \end{array} \\
 N \times d_k \qquad \qquad \qquad N \times N \qquad \qquad \qquad N \times d_v \qquad \qquad \qquad N \times d_v
 \end{array}$$

$$\begin{array}{c}
 \left[ \begin{array}{c} Q \\ \vdots \\ q4 \end{array} \right] \times \begin{array}{c} K^T \\ \vdots \\ \text{k1} | \text{k2} | \text{k3} | \text{k4} \end{array} = \begin{array}{c} QK^T \\ \vdots \\ \text{q4}\cdot\text{k1} | \text{q4}\cdot\text{k2} | \text{q4}\cdot\text{k3} | \text{q4}\cdot\text{k4} \end{array} \right] \times \begin{array}{c} V \\ \vdots \\ \text{v1} \\ \text{v2} \\ \text{v3} \\ \text{v4} \end{array} = \begin{array}{c} A \\ \vdots \\ \text{a4} \end{array} \\
 1 \times d_k \qquad \qquad \qquad 1 \times N \qquad \qquad \qquad N \times d_v \qquad \qquad \qquad 1 \times d_v
 \end{array}$$

# Inference Math

- Decoding is memory-bandwidth bound: Need to load Params + KV cache to memory
- Model bandwidth utilization (MBU) = No. of bytes loaded per sec / Theoretical mem bw
- Typical MBU: 30-70%

gpt-fast, A100 (max 2TB/s):

Model	Technique	Tokens/Second	Memory Bandwidth (GB/s)
Llama-2-7B	Base	104.9	1397.31
	8-bit	155.58	1069.20
	4-bit (G=32)	196.80	862.69
Llama-2-70B	Base	OOM	
	8-bit	19.13	1322.58
	4-bit (G=32)	25.25	1097.66
Llama-3.1-8B	Base	93.89	1410.76
	8-bit	137.64	1030.89
Llama-3.1-70B	Base	OOM	
	8-bit	18.04	1253.78

gpt-fast: <https://github.com/pytorch-labs/gpt-fast>

# Inference: Reducing KV cache with Grouped Query Attention

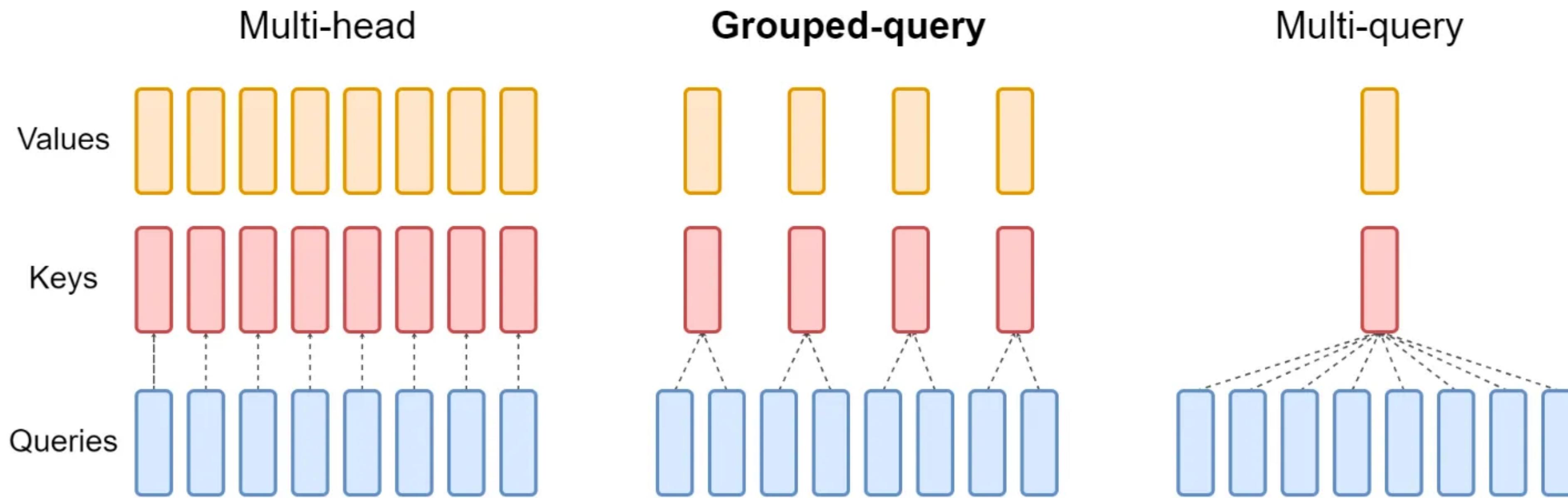
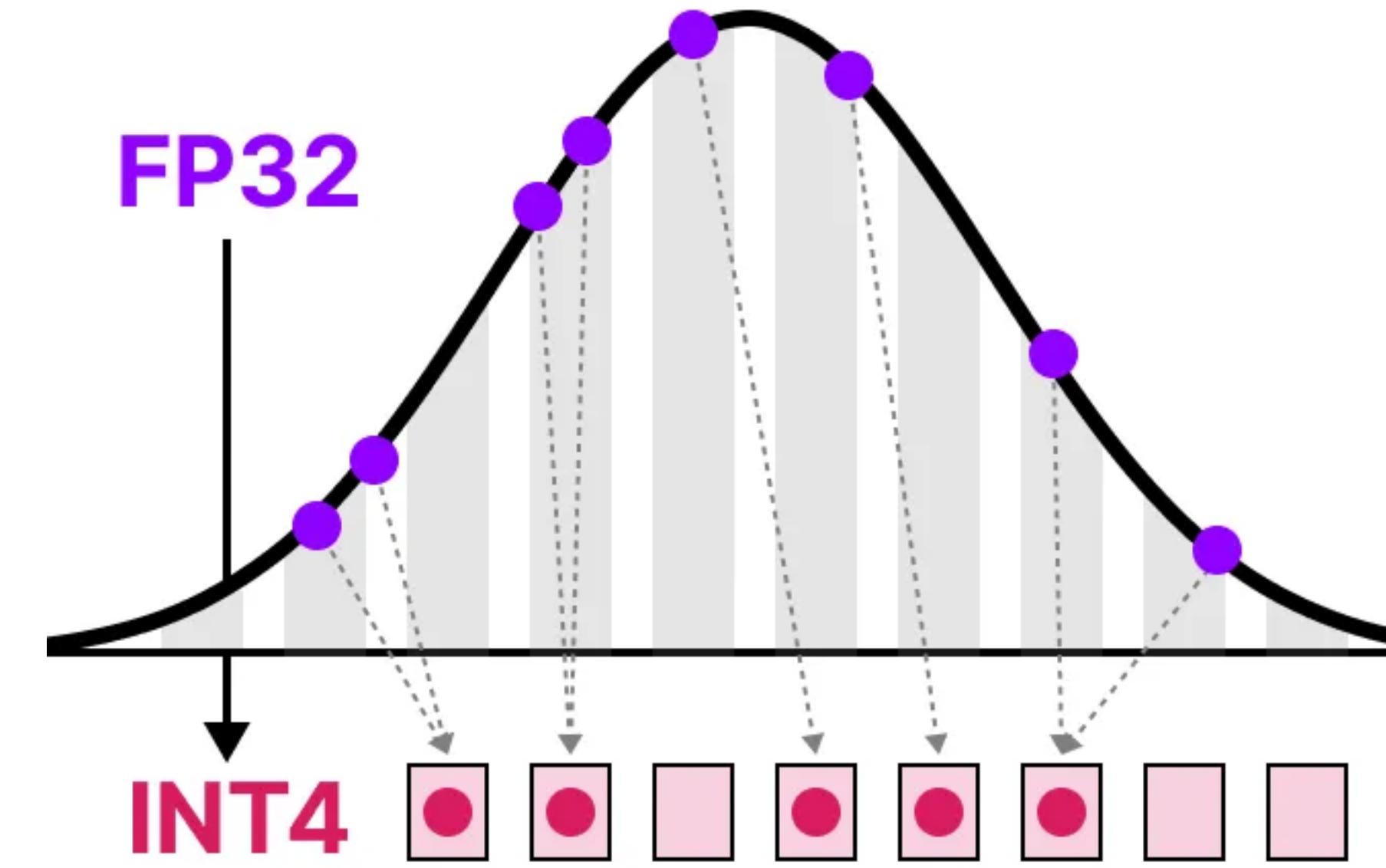


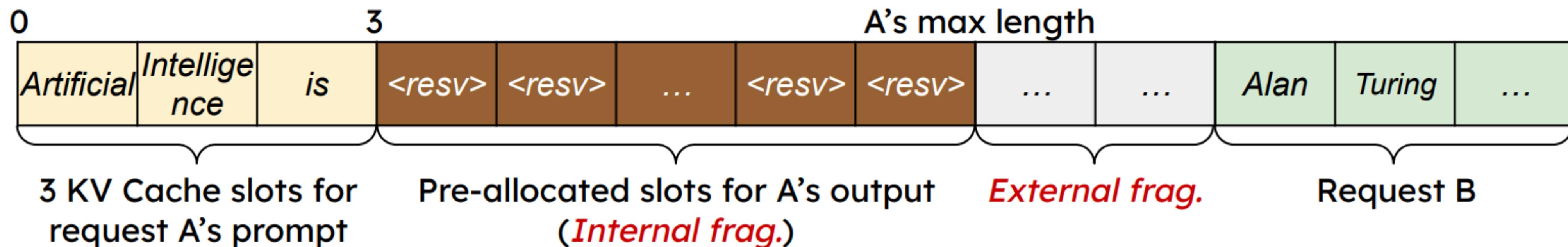
Figure 2: Overview of grouped-query method. Multi-head attention has  $H$  query, key, and value heads. Multi-query attention shares single key and value heads across all query heads. Grouped-query attention instead shares single key and value heads for each *group* of query heads, interpolating between multi-head and multi-query attention.

# Inference: Weight & KV Cache Quantization



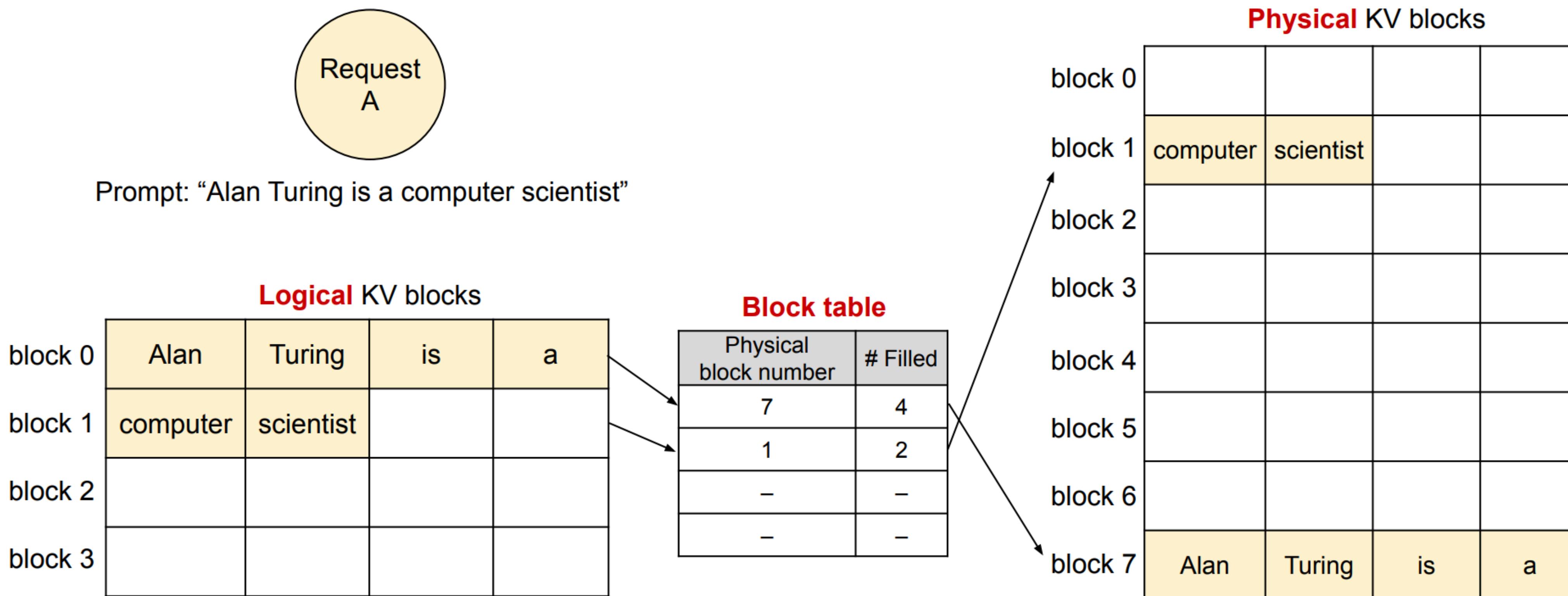
- Reduce no. bytes loaded, at the cost of slightly lower model quality

# Inference: Challenge with KV cache memory management



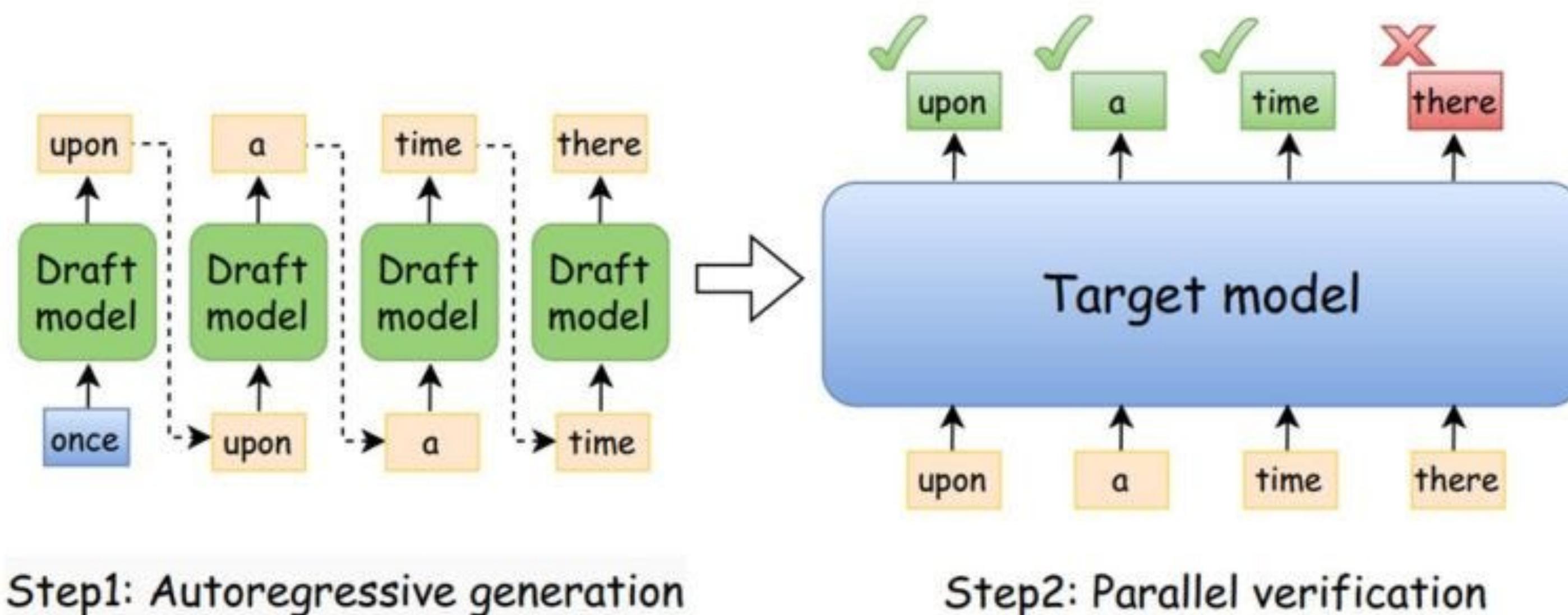
- **Pre-allocates contiguous** memory to the request's max length
- **Memory fragmentation:**
  - Internal fragmentation due to unknown output length
  - External fragmentation due to non-uniform per-request max length

# Inference: Memory management with Paged KV

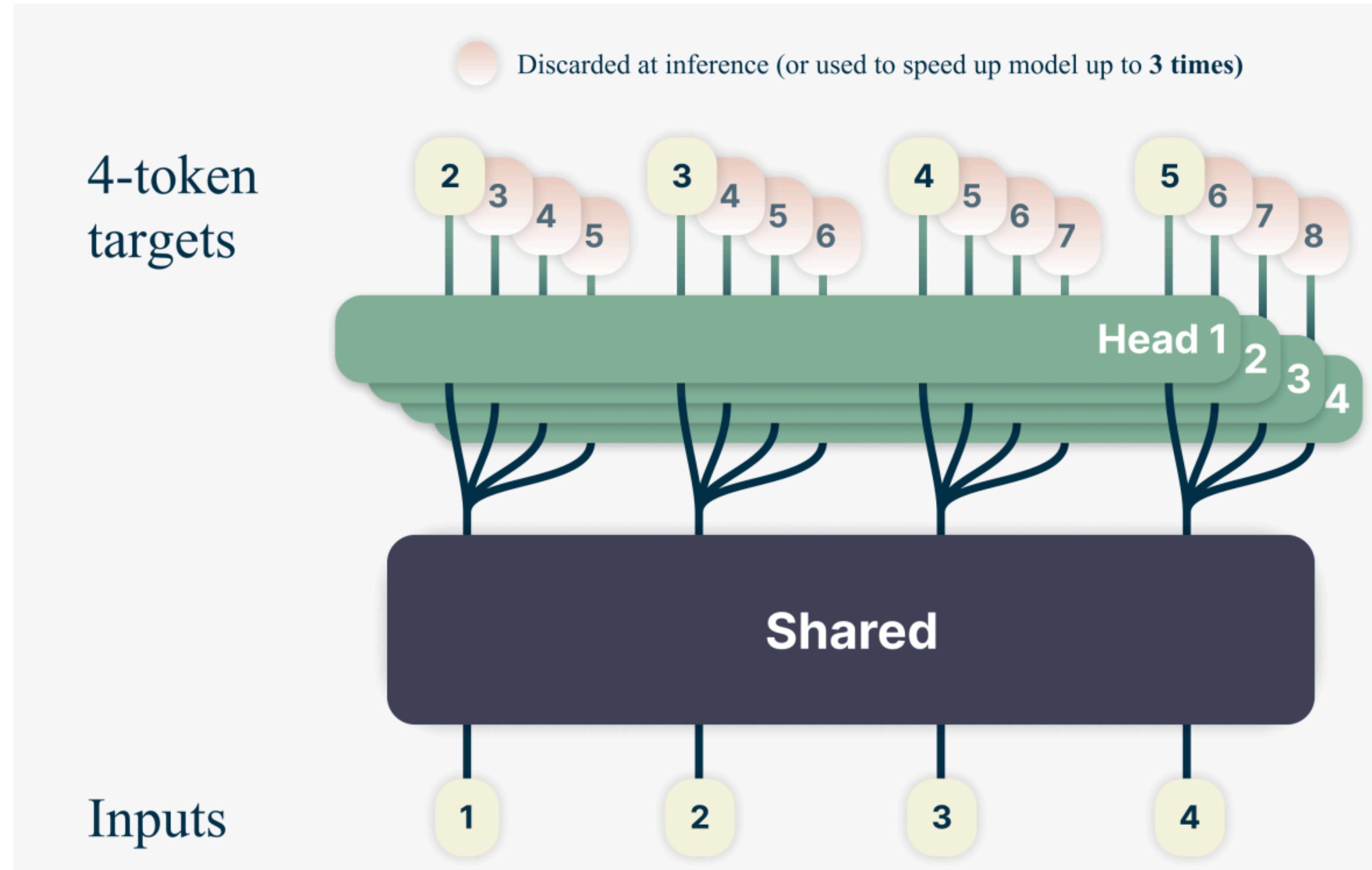


# Inference: Speculative Decoding

- Use a small Draft model to predict what the Target model will output
- Verify with Target model multiple tokens in parallel ( $\approx$  same cost as processing 1 token)



# Multi-token Prediction



# Concluding Thoughts

- Efficient training & inference for LLM is a wide open field
- Intersection of model architecture, systems, and algorithms
- Trend: close co-design of model, inference system, and hardware