

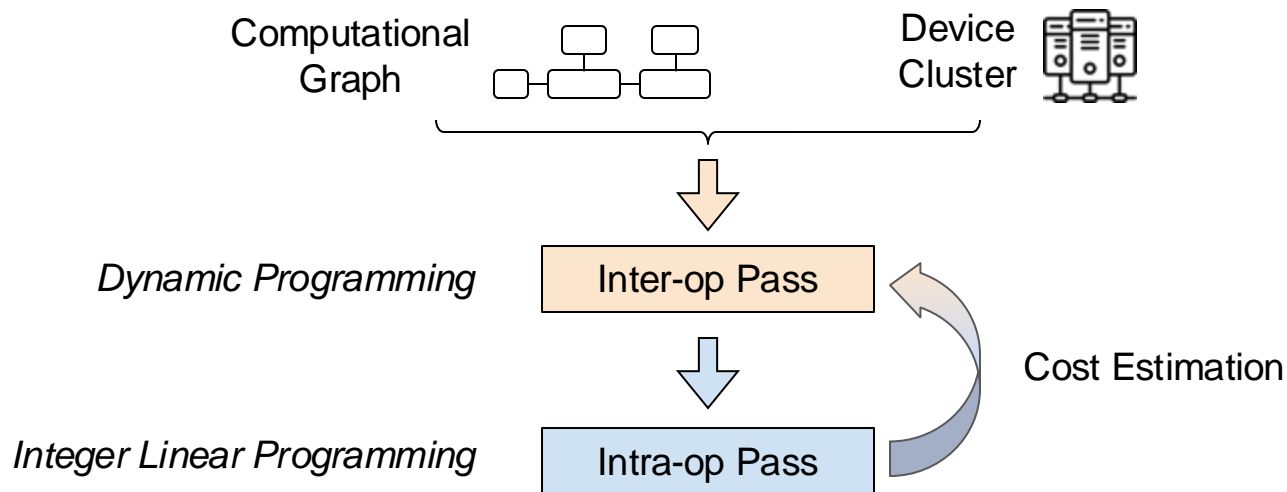
# Logistics

- PA3 is posted.
  - Two programming assignments
    - One MoE
    - One LLM inference
  - One theoretical assignment
    - Scaling law
  - One essay

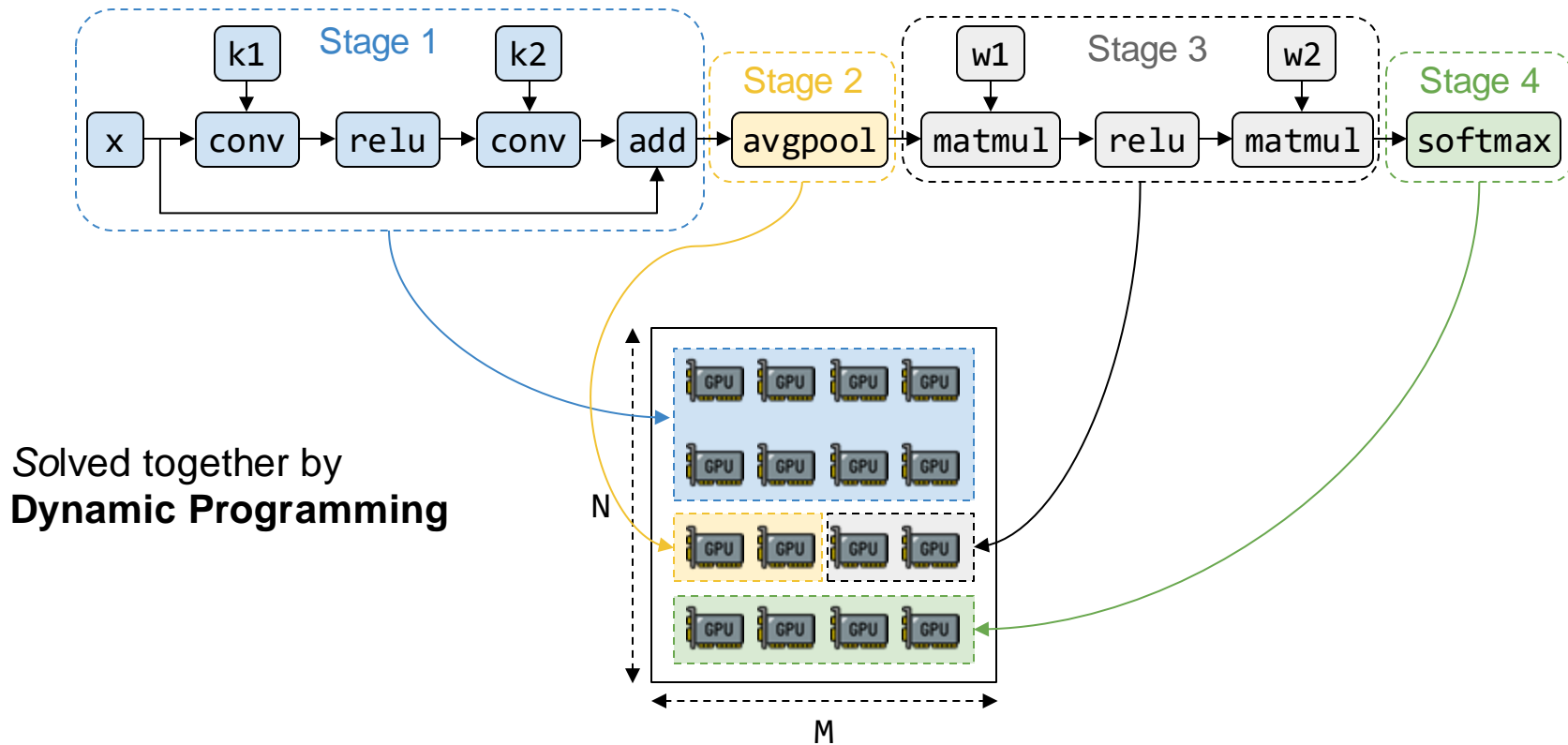
# Where We Are

- Motivation
- History
- Parallelism Overview
- Data parallelism
- Model parallelism
  - Inter-op parallelism
  - Intra-op parallelism
- **Auto-parallelization**

# Alpa Compiler: Hierarchical Optimization

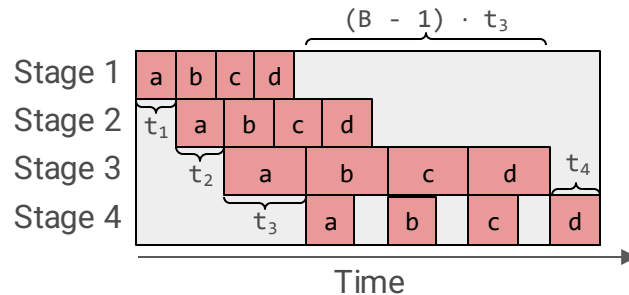
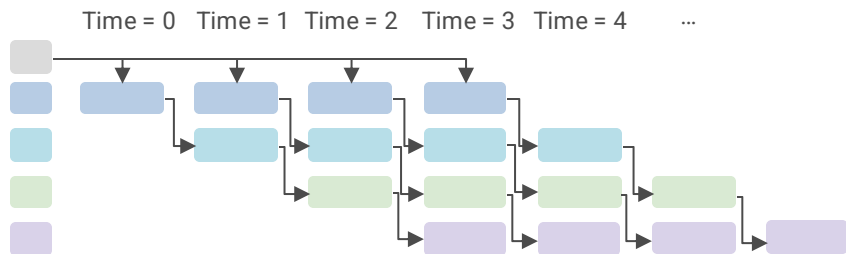


# Inter-op Pass



# Pipeline Execution Latency

Inter-op Pass



warmup phase

stable phase

$$T = \sum_i^S t_i + (B - 1) \cdot \max_{1 \leq j \leq S} \{t_j\}$$

# Inter-op Pass: Dynamic Programming

**Optimization objective:** Find the optimal (stage, mesh) pairs that minimize  $T$ .

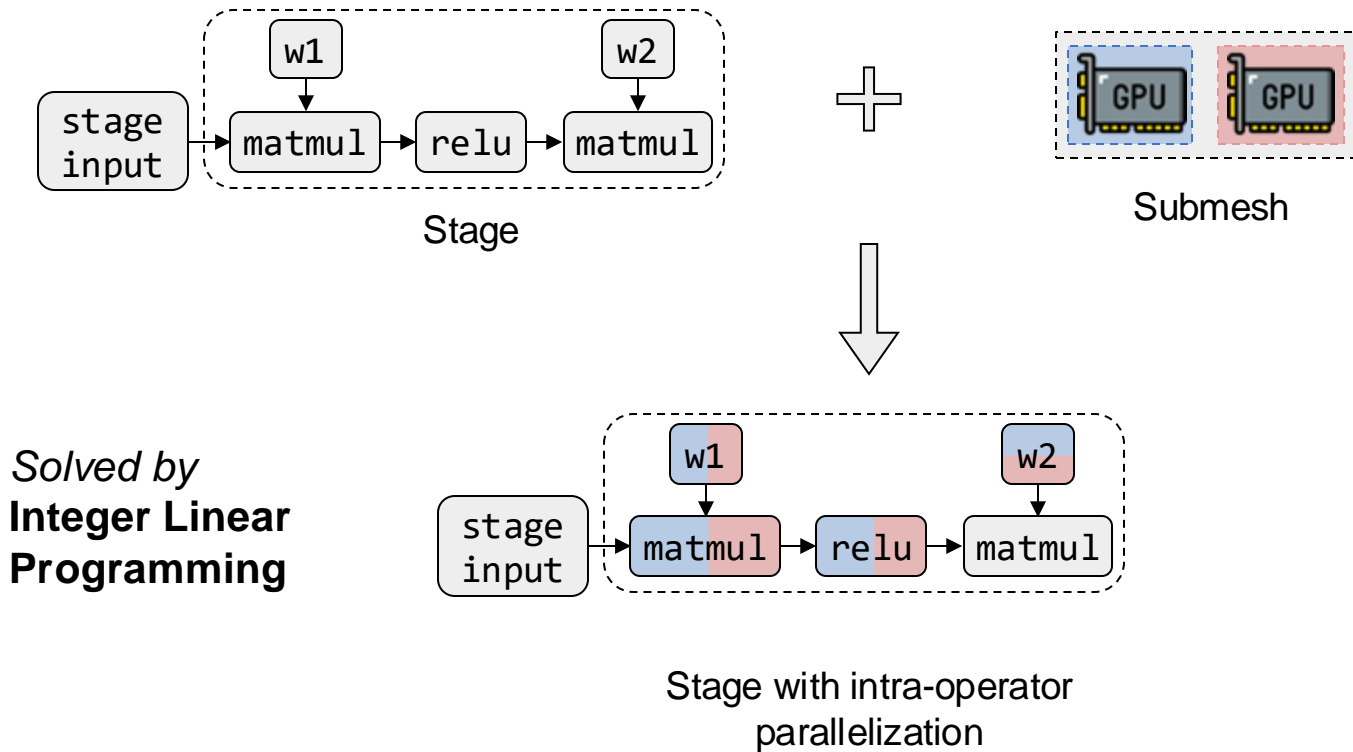
$$T = \underbrace{\sum_{i=1}^S t_i}_{\text{warmup phase}} + \underbrace{(B-1) \cdot \max_{1 \leq j \leq S} \{t_j\}}_{\text{stable phase}}$$

the **optimal** latency of executing stage  $i$  on its assigned mesh  $i$  :  $t_i = t_{\text{intra}}^*(\text{stage}_i, \text{mesh}_i)$

**Solution:**

Enumerate all possible  $\max_{1 \leq j \leq S} \{t_j\}$  (stable phase) and convert the first term  $\sum_{i=1}^S t_i$  (warmup phase) into a 2-dimensional knapsack problem.

# Intra-op Pass

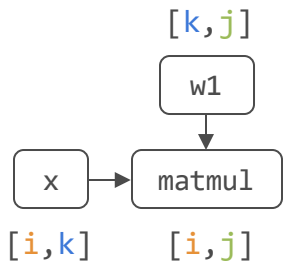


# Intra-op Pass: Computation

Intra-op Pass

Device 1  
Device 2

Row-partitioned Column-partitioned Replicated



$$\text{matmul}[i, j] = \sum_k x[i, k] \times w1[k, j] \quad \text{Cost}$$

Algo#1: loop  $i$   =   $\times$   Cost1

Algo#2: loop  $j$   =   $\times$   Cost2

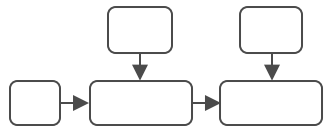
Algo#3: loop  $k$    =   $\times$   Cost3

Algo#4: ...



# Intra-op Pass: Communication

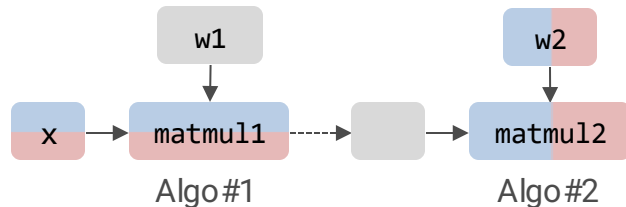
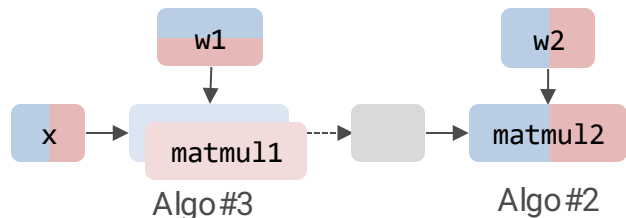
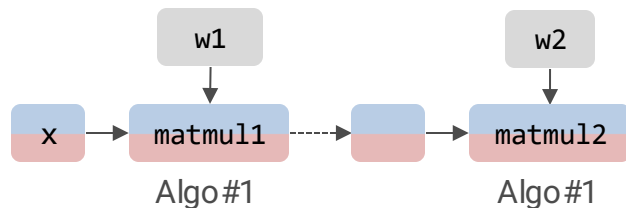
Intra-op Pass



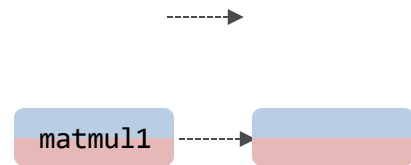
Algo#1:  =  × 

Algo#2:  =  × 

Algo#3:   =  × 

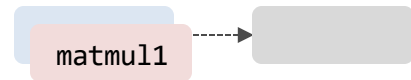


Layout Conversion



Cost

0



all-reduce



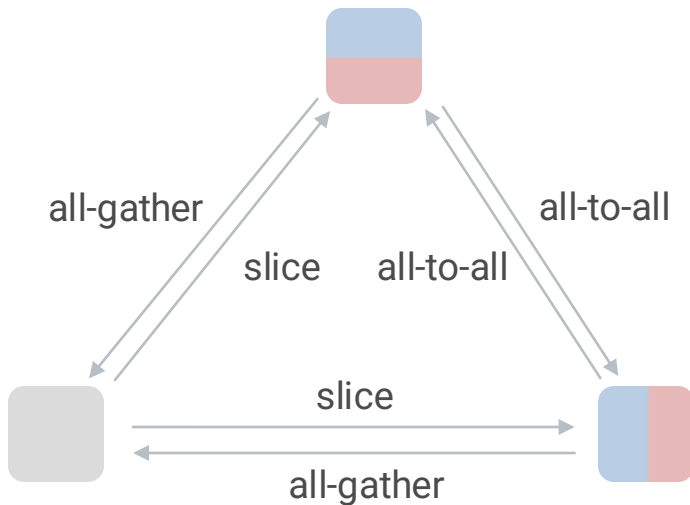
all-gather

# Intra-op Pass: Layout Conversion

Intra-op Pass

Device 1  
Device 2

Row-partitioned Column-partitioned Replicated



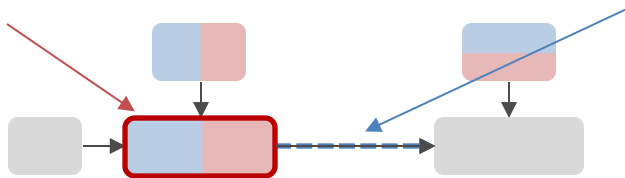
# Intra-op Pass: ILP Formulation

Intra-op Pass

Goal: Within each stage, “color” every node in the stage, so the execution latency of this stage on its assigned mesh is minimized.

For every node (op), enumerate all possible parallel algorithms

For every edge, infer the cost due to layout conversion

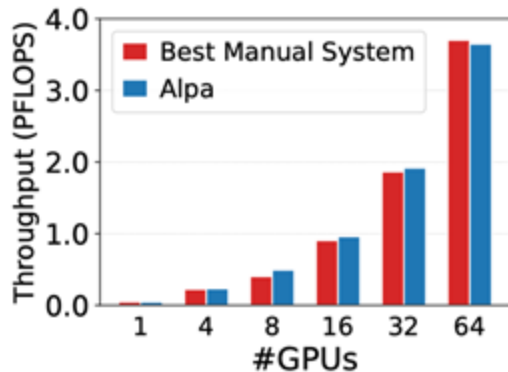


Minimize **node-cost** + **edge-cost**

s.t. peak memory usage < memory budget

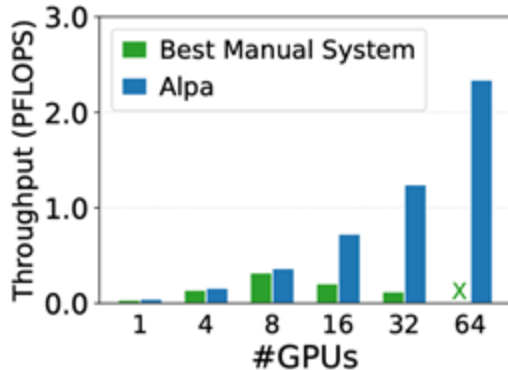
# Evaluation: Comparing with Previous Works

## GPT (up to 39B)



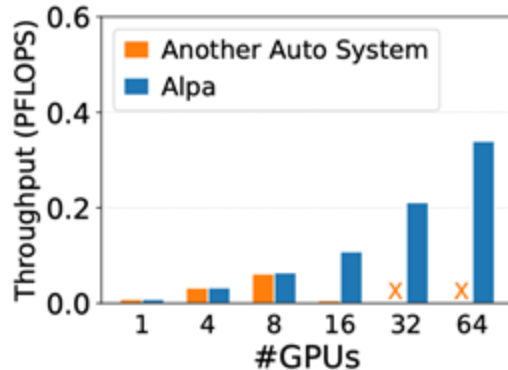
Match specialized manual systems.

## GShard MoE (up to 70B)



Outperform the manual baseline by up to 8x.

## Wide-ResNet (up to 13B)



Generalize to models without manual plans.

*Weak scaling results where the model size grow with #GPUs.*

*Evaluated on 8 AWS EC2 p3.16xlarge nodes with 8 16GB V100s each (64 GPUs in total).*

# Automatic Parallelization Methods

## Search-based methods

- ✓ Easy to extend the search space
- ✓ No training cost
- ✗ High inference cost
- ✗ Not explainable
- ✗ No optimality guarantee

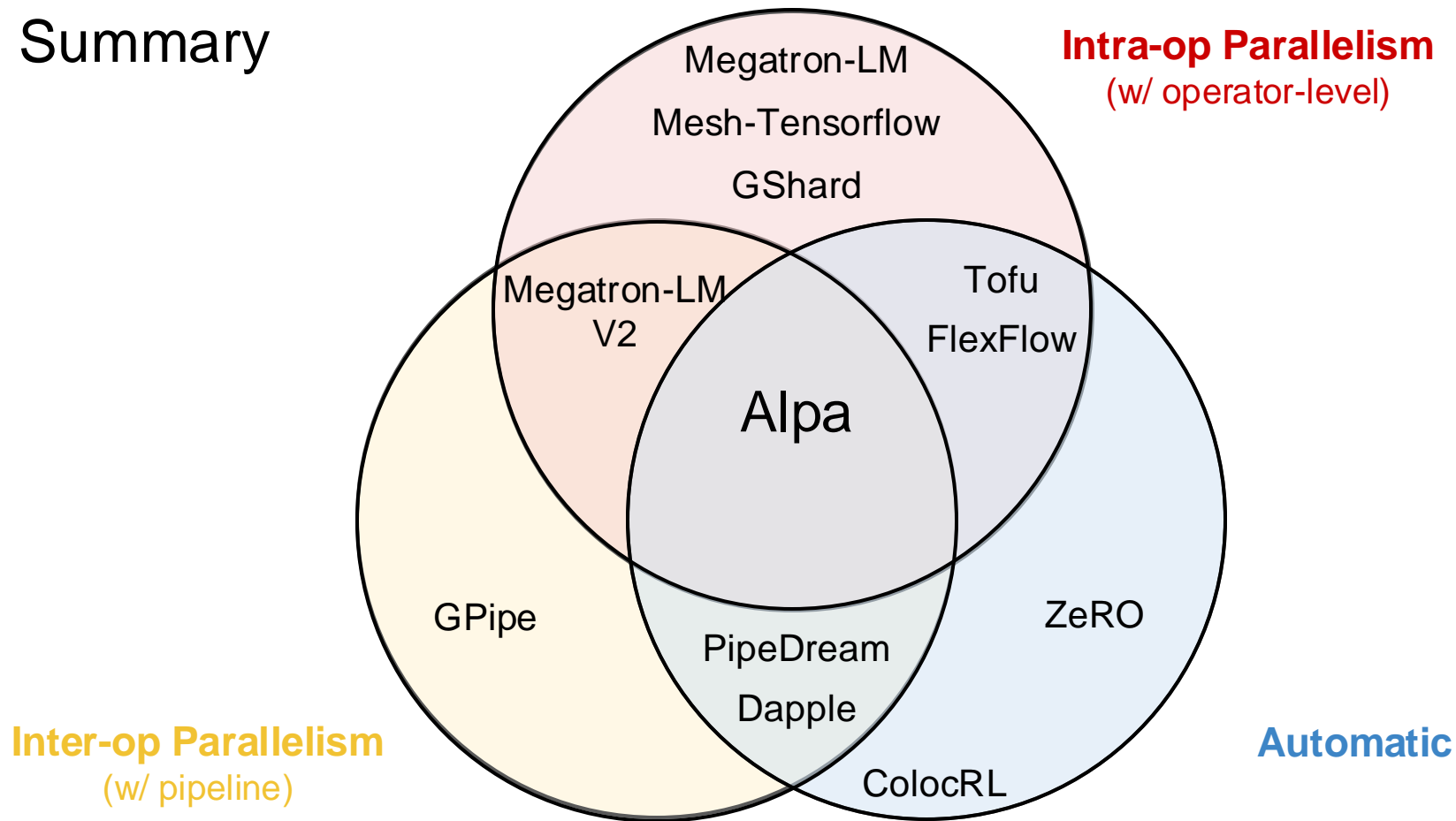
## Learning-based methods

- ✓ Easy to extend the search space
- ✗ High training cost
- ✓ Low inference cost
- ✗ Not explainable
- ✗ No optimality guarantee

## Optimization-based methods

- ✗ Non-trivial to extend the search space
- ✓ No training cost
- ✓ Medium inference cost
- ✓ Explainable
- ✓ Some optimality guarantee

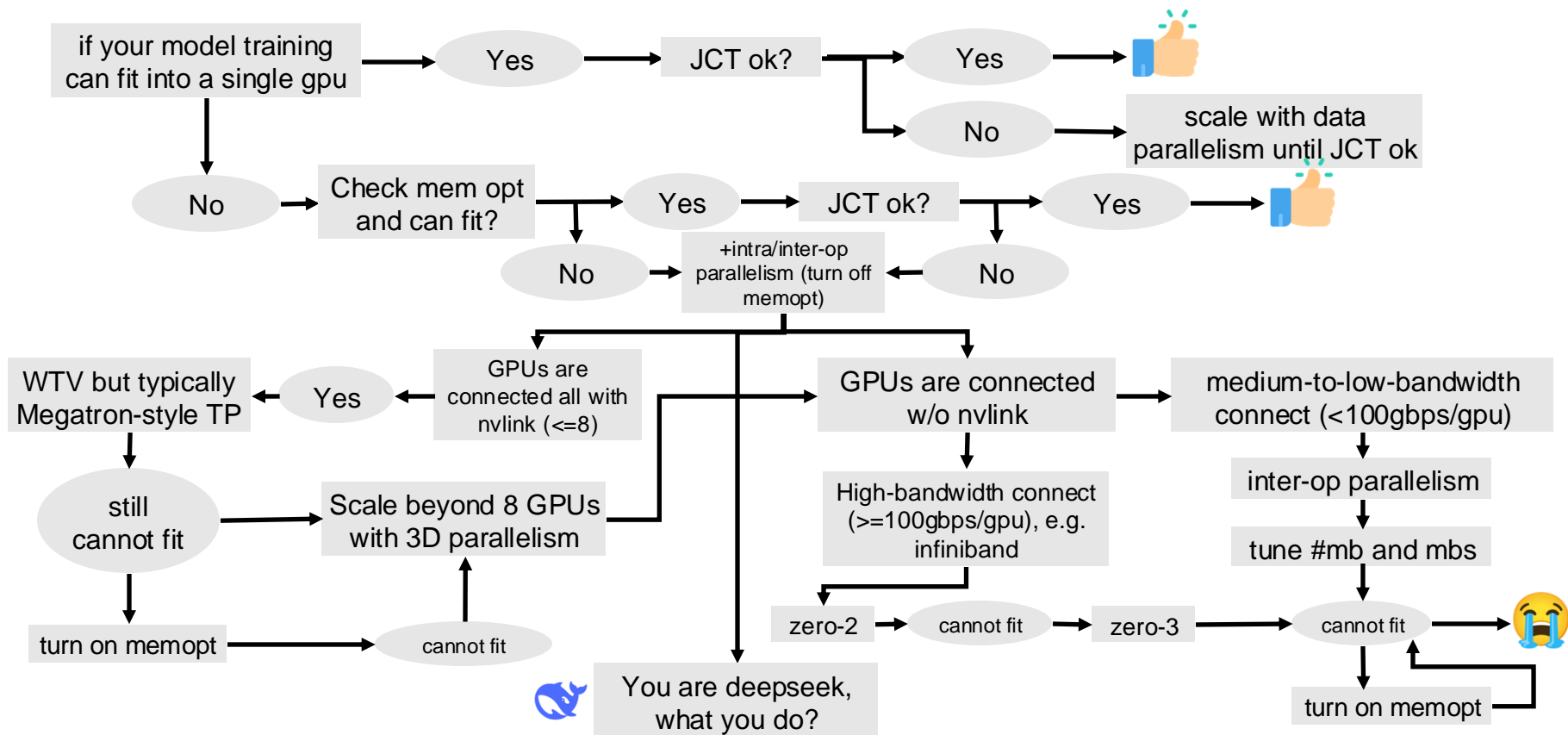
# Summary



# Summary: How to Choose Parallelism

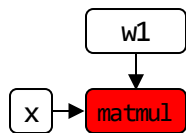
1. Use automatic compiler if not transformer
2. Manual parallelism search for transformers:
  - Factors to consider
    - #GPUs you have
    - Model size
    - JCT (Job completion time)
    - Communication bandwidth
    - etc.

# Hao's Ultimate Guide





# Big Picture



Dataflow Graph

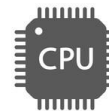
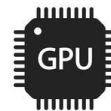
Autodiff

Graph Optimization

Parallelization

Runtime: schedule / memory

Operator optimization/compilation



## Next: Connecting the Dots

LLMSys

Optimizations and Parallelization

MLSys Basics

# Large Language Models

- Transformers, Attentions
- Scaling Law
  - MoE
- Connecting the dots: Training Optimizations
  - Flash attention
  - Long context, parallelism
- Serving and inference optimization
  - Continuous batching and Paged attention
  - Speculative decoding (Guest Lecture)
- Connecting the dots: Deepseek-v3
- Hot topics

## Next Token Prediction

$$P(\textit{next word} \mid \textit{prefix})$$

San Diego has very nice _	surfing	0.4
	weather	0.5
	snow	0.01
San Francisco is a city of _	innovation	0.6
	homeless	0.3

## Next Token Prediction

$$\begin{aligned} & \text{Probability("San Diego has very nice weather")} \\ &= P(\text{"San Diego"}) P(\text{"has"} | \text{"San Diego"}) P(\text{"very"} | \text{"San Diego has"}) P(\text{"city"} | \dots) \dots P(\text{"weather"} | \dots) \end{aligned}$$

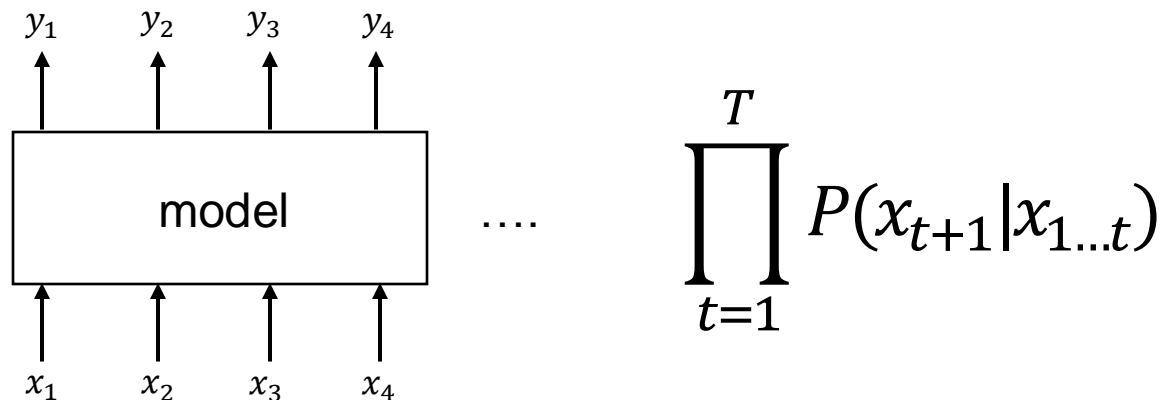
$$\text{Max Prob}(x_{1:T}) = \prod_{t=1}^T P(x_{t+1} | x_{1..t})$$

MLE on observed data  $x_{1:T}$ ,

This is next token prediction.  
Predicting using seq2seq NNs.

# Sequence Prediction

Take a set of input sequence, predict the output sequence



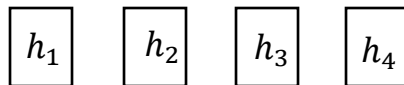
Predict each output based on history  $y_t = f_{\theta}(x_{1:t})$

There are many ways to build up the predictive model

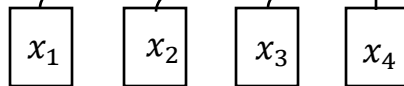
# “Attention” Mechanism

Generally refers to the approach that weighted combine individual states

Attention output



Hidden states from previous layer



$$h_t = \sum_{i=1}^t s_i x_t$$

Intuitively  $s_i$  is “attention score” that computes how relevant the position  $i$ ’s input is to this current hidden output

There are different methods to decide how attention score is being computed

# Self-Attention Operation

Self attention refers to a particular form of attention mechanism.

Given three inputs  $Q, K, V \in \mathbb{R}^{T \times d}$  (“queries”, “keys”, “values”)

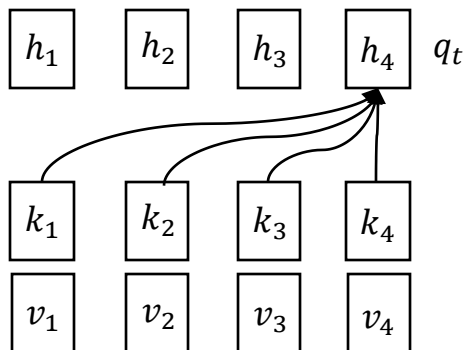
Define the self-attention as:

$$\text{SelfAttention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{d^{1/2}}\right) V$$



# A Closer Look at Self-Attention

Use  $q_t, k_t, v_t$  to refers to row  $t$  of the  $K$  matrix



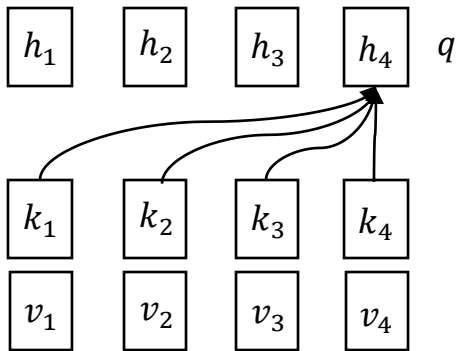
Ask the following question:

How to compute the output  $h_t$ , based on  $q_t, K, V$  one timestep  $t$

To keep presentation simple, we will drop suffix  $t$  and just use  $q$  to refer to  $q_t$  in next few slide

## A Closer Look at Self-Attention

Use  $q_t, k_t, v_t$  to refers to row  $t$  of the  $K$  matrix



Conceptually, we compute the output in the following two steps:

Pre-softmax “attention score”

$$s_i = \frac{1}{\sqrt{d}} q k_i^T$$

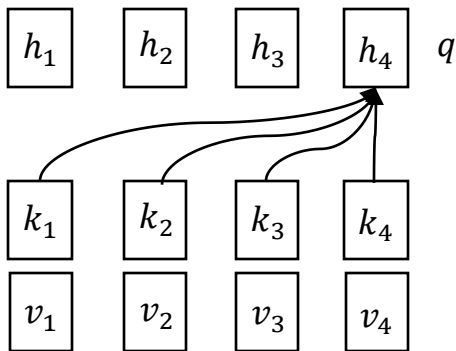
Weighted average via softmax

$$h = \sum_i \text{softmax}(s)_i v_i = \frac{\sum_i \exp(s_i) v_i}{\sum_j \exp(s_j)}$$

Intuition:  $s_i$  computes the relevance of  $k_i$  to the query  $q$ , then we do weighted sum of values proportional to their relevance

# Comparing the Matrix Form and the Decomposed Form

Use  $q_t, k_t, v_t$  to refers to row  $t$  of the  $K$  matrix



$$\text{SelfAttention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{d^{1/2}}\right) V$$

Pre-softmax “attention score”

$$S_{ti} = \frac{1}{\sqrt{d}} q_t k_i^T$$

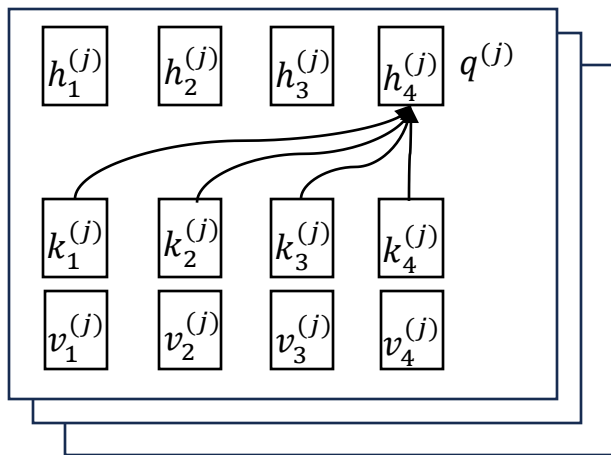
Weighed average via softmax

$$h_t = \sum_i \text{softmax}(S_{t,:})_i v_i = \text{softmax}(S_{t,:}) V$$

Intuition:  $s_i$  computes the relevance of  $k_i$  to the query  $q$ ,  
then we do weighted sum of values proportional to their relevance

# Multi-Head Attention

Have multiple “attention heads”  $Q^{(j)}, K^{(j)}, V^{(j)}$  denotes  $j$ -th attention head



Apply self-attention in each attention head

$$\text{SelfAttention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{d^{1/2}}\right) V$$

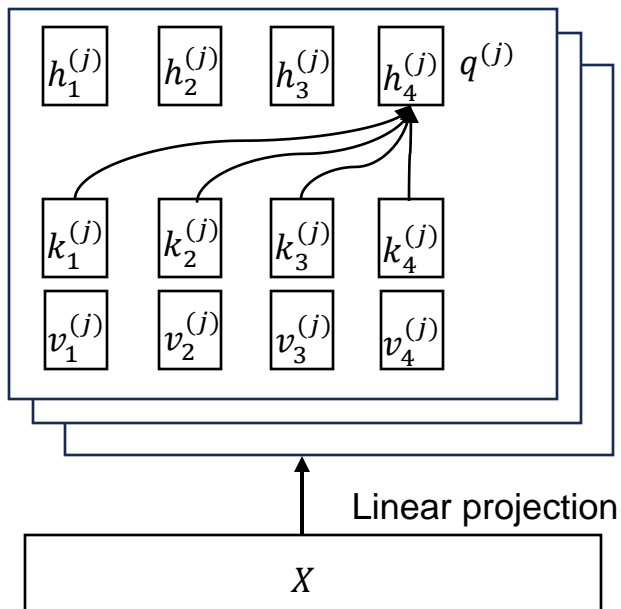
Concatenate all output heads together as output

Each head can correspond to different kind of information.

Sometimes we can share the heads: GQA(group query attention) all heads share  $K, V$  but have different  $Q$

## How to get Q K V?

Obtain  $Q, K, V$  from previous layer's hidden state  $X$  by linear projection



$$Q = XW_q$$

$$K = XW_K$$

$$V = XW_V$$

Can compute all heads and  $Q, K, V$  together then split/reshape out into individual  $Q, K, V$  with multiple heads

# Transformer Block

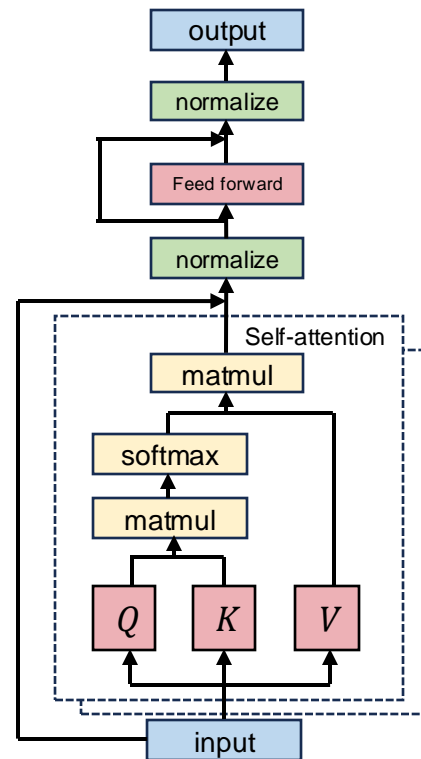
A typical transformer block

$$Z = \text{SelfAttention}(XW_K, XW_Q, XW_V)$$

$$Z = \text{LayerNorm}(X + Z)$$

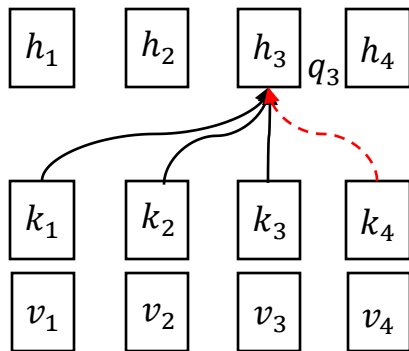
$$H = \text{LayerNorm}(\text{ReLU}(ZW_1)W_2 + Z)$$

(multi-head) self-attention, followed by a linear layer and ReLU and some additional residual connections and normalization



# Masked Self-Attention

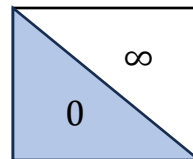
In the matrix form, we are computing weighted average over all inputs



In auto regressive models, usually it is good to maintain casual relation, and only attend to some of the inputs (e.g. skip the red dashed edge on the left). We can add “attention mask”

$$\text{MaskedSelfAttention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{d^{1/2}} - M\right)V$$

$$M_{ij} = \begin{cases} \infty, & j > i \\ 0, & j \leq i \end{cases}$$

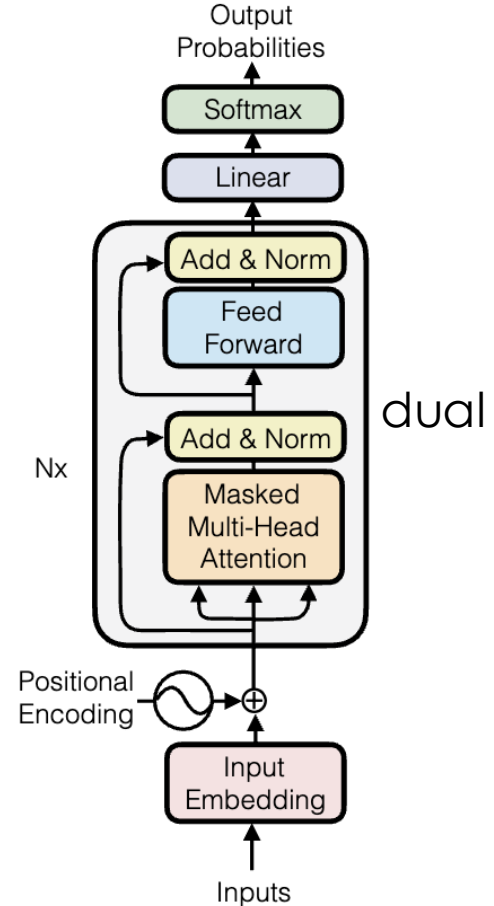


Only attend to previous inputs. Depending on input structure and model, attention mask can change.

We can also simply skip the computation that are masked out if there is a special implementation to do so

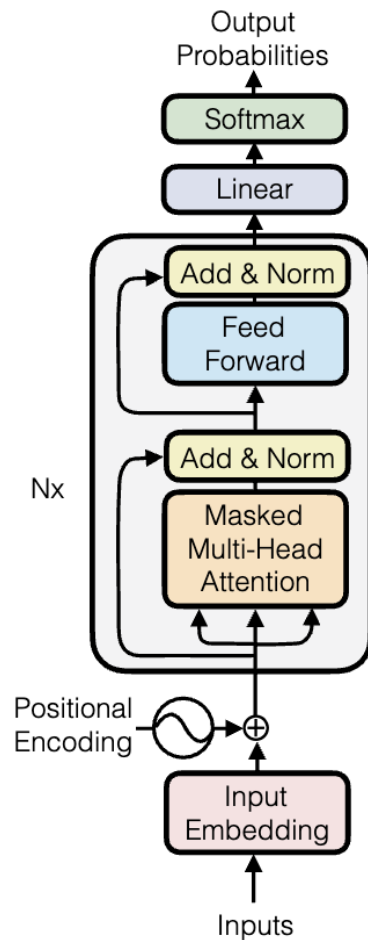
# Transformers

- Transformer decoders
  - Many of them
  - Really just: attentions + layernorm + MLPs
- Word embeddings
- Position embeddings
  - Rotary embedding
- Loss function: cross entropy loss over a sequence



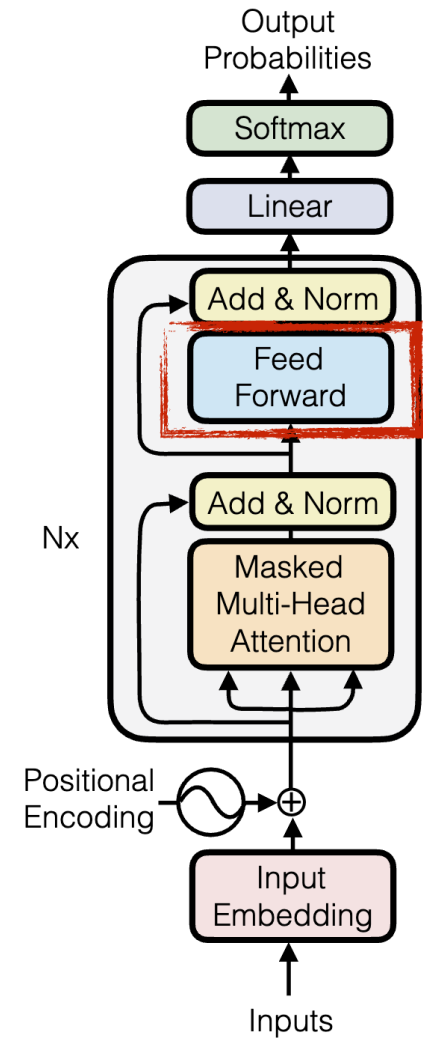
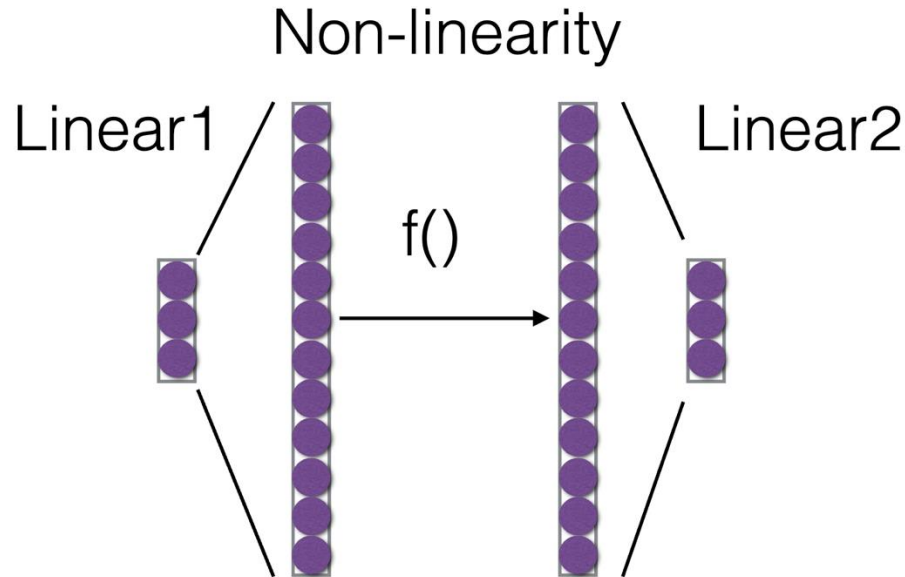


# Transformers



# Feedforward Layers

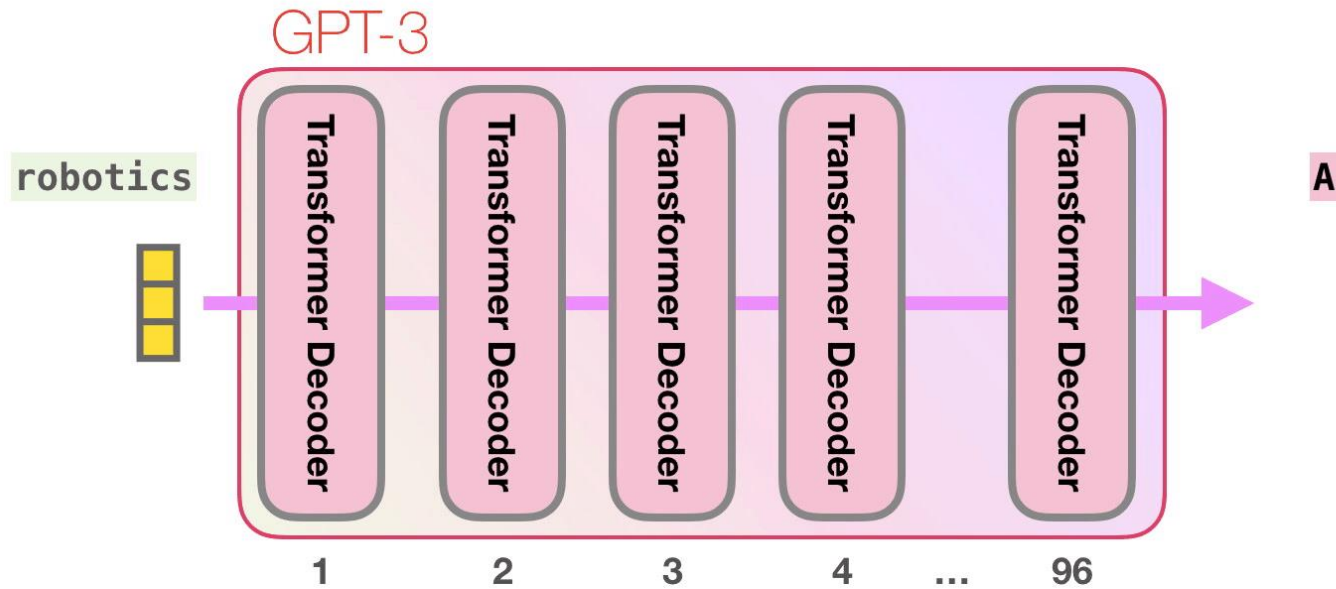
$$\text{FFN}(x; W_1, \mathbf{b}_1, W_2, \mathbf{b}_2) = f(\mathbf{x}W_1 + \mathbf{b}_1)W_2 + \mathbf{b}_2$$



# Computing Components in LLMs?

- Transformer decoders (many of them)
  - self-attentions (slow)
  - layernorm, residual (fast)
  - MLPs (slow)
  - Nonlinear (fast)
- Word embeddings (fast)
- Position embeddings (fast)
  - Absolute embedding vs. relative embedding
- Loss function: cross entropy loss over a sequence of words

# LLMs

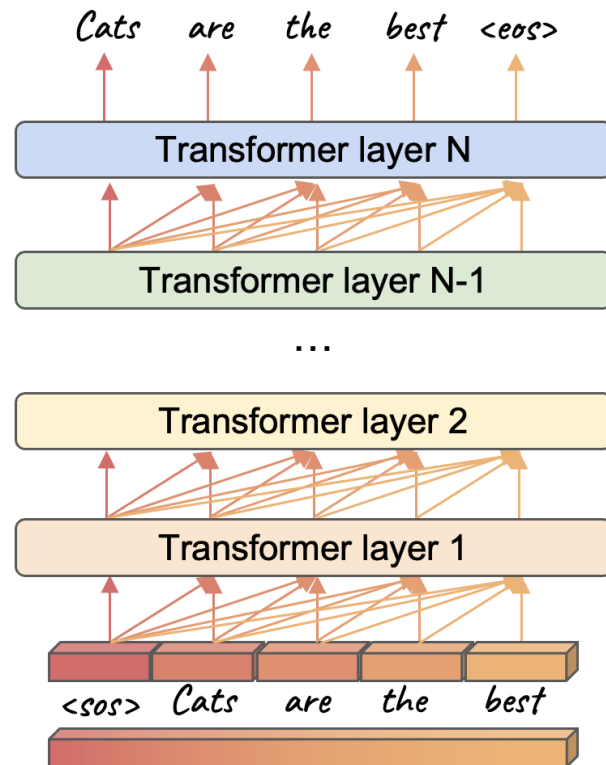


# Original Transformer vs. LLM today

	Vaswani et al.	LLaMA
<b>Norm Position</b>	Post	Pre
<b>Norm Type</b>	LayerNorm	RMSNorm
<b>Non-linearity</b>	ReLU	SiLU
<b>Positional Encoding</b>	Sinusoidal	RoPE

# Training LLMs

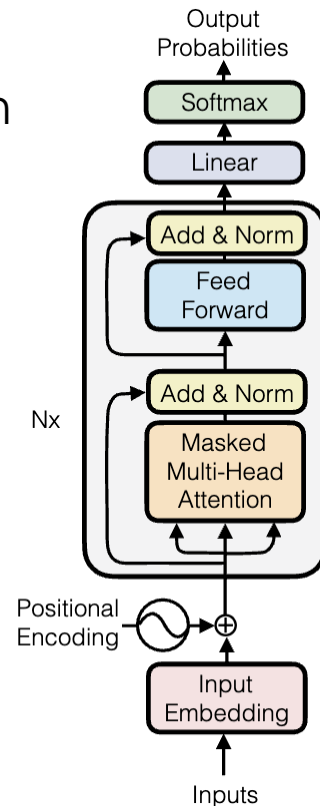
- Sequences are **known a priori**
- For each position, look at  $[1, 2, \dots, t-1]$  words to predict word  $t$ , and calculate the loss at  $t$
- Parallelize the computation across all token positions, and then apply masking



# Connecting the Dots: Compute/Comm characteristic of LLMs

Key characteristics: compute, memory, communication

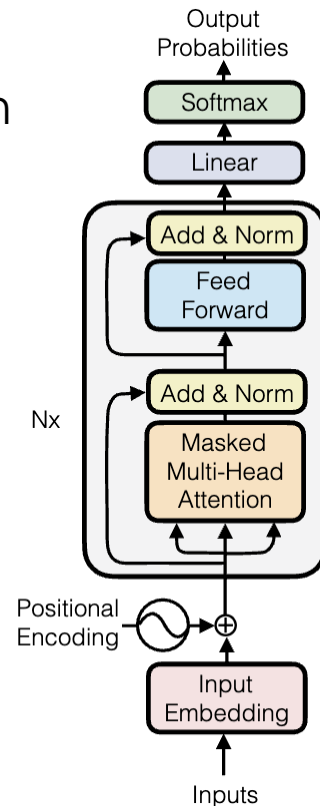
- calculate the number of parameters of an LLM?
  - memory, communication
- calculate the flops needed to train an LLM?
  - compute
- calculate the memory needed to train an LLM?
  - memory, communication



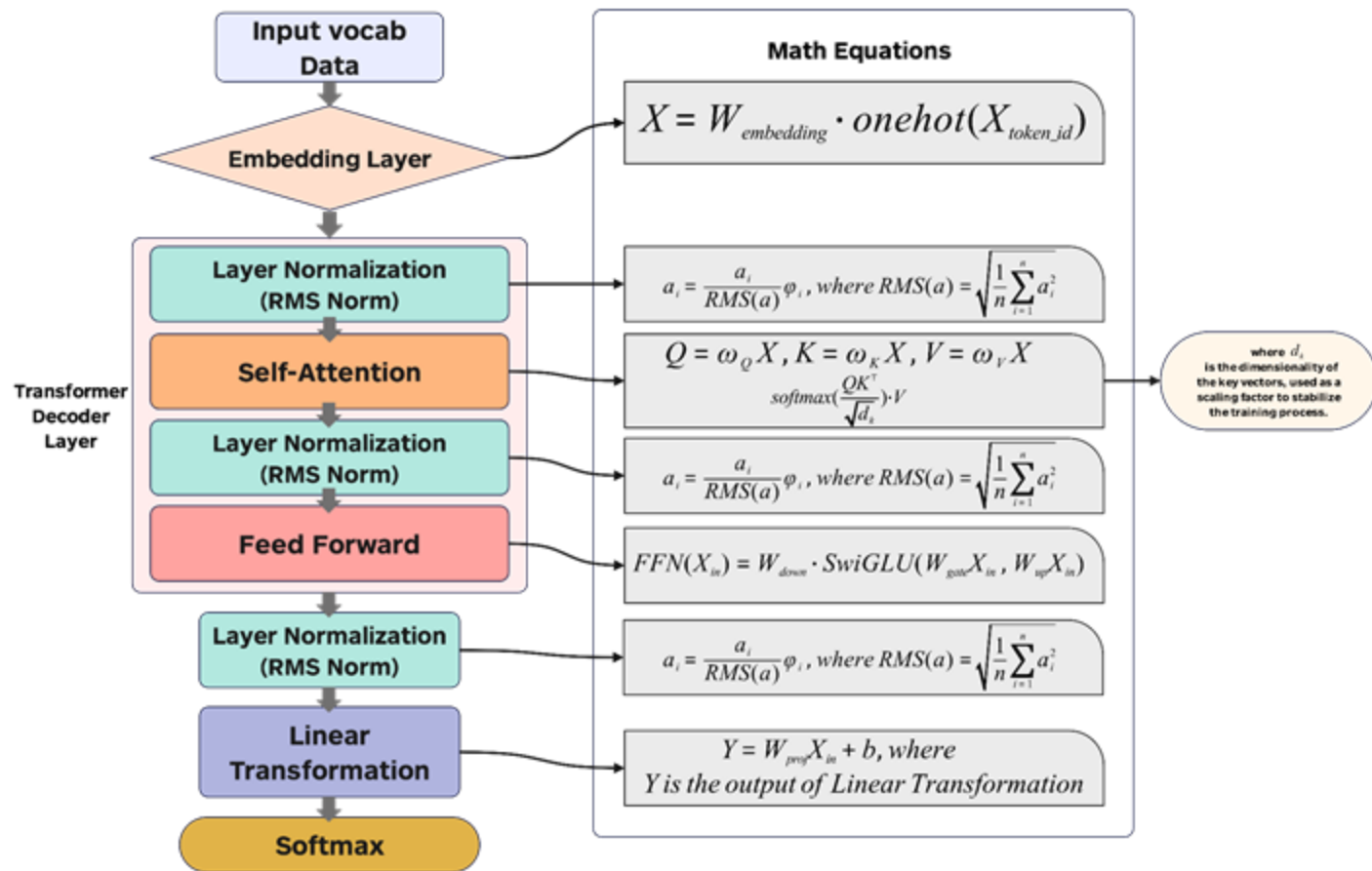
# Connecting the Dots: Compute/Comm characteristic of LLMs

Key characteristics: compute, memory, communication

- calculate the number of parameters of an LLM?
- calculate the flops needed to train an LLM?
- calculate the memory needed to train an LLM?







# Feed Forward SwiGLU

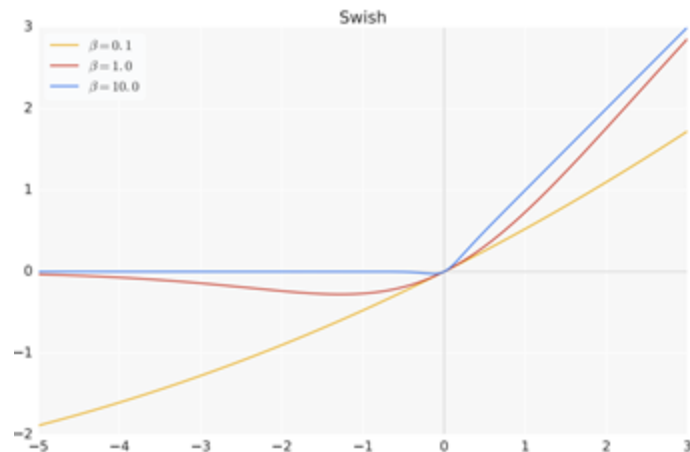
The general formula for SwiGLU is:

$$\text{SwiGLU}(x) = \text{Swish}(xW_1 + b_1) \odot (xW_2 + b_2)$$

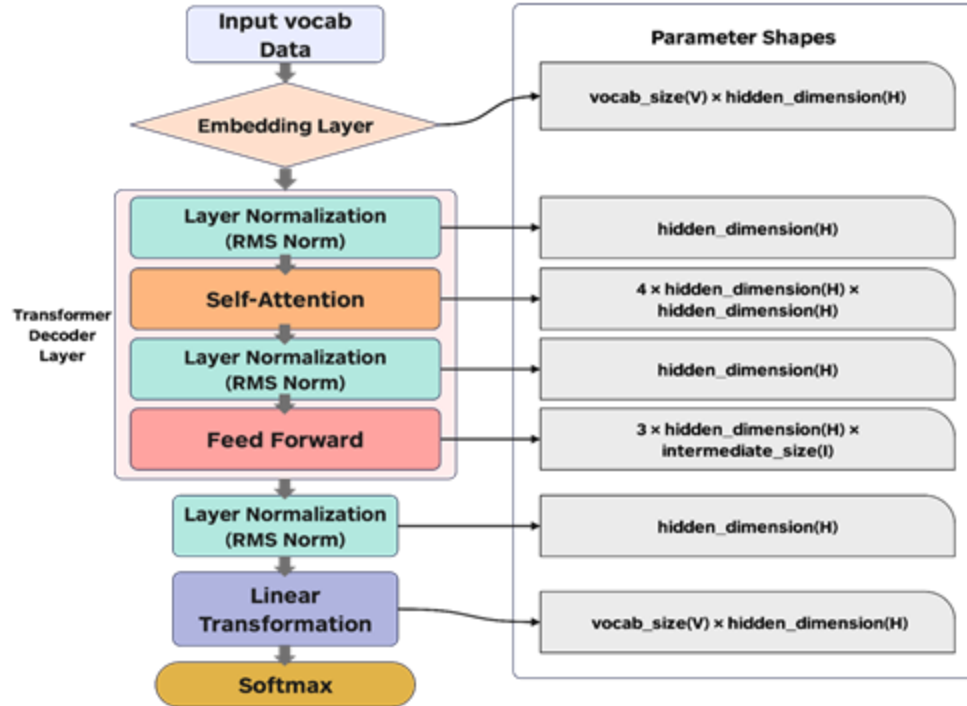
**Swish** is the activation function applied to one branch, defined as:

$$\text{Swish}(z) = z \cdot \sigma(z)$$

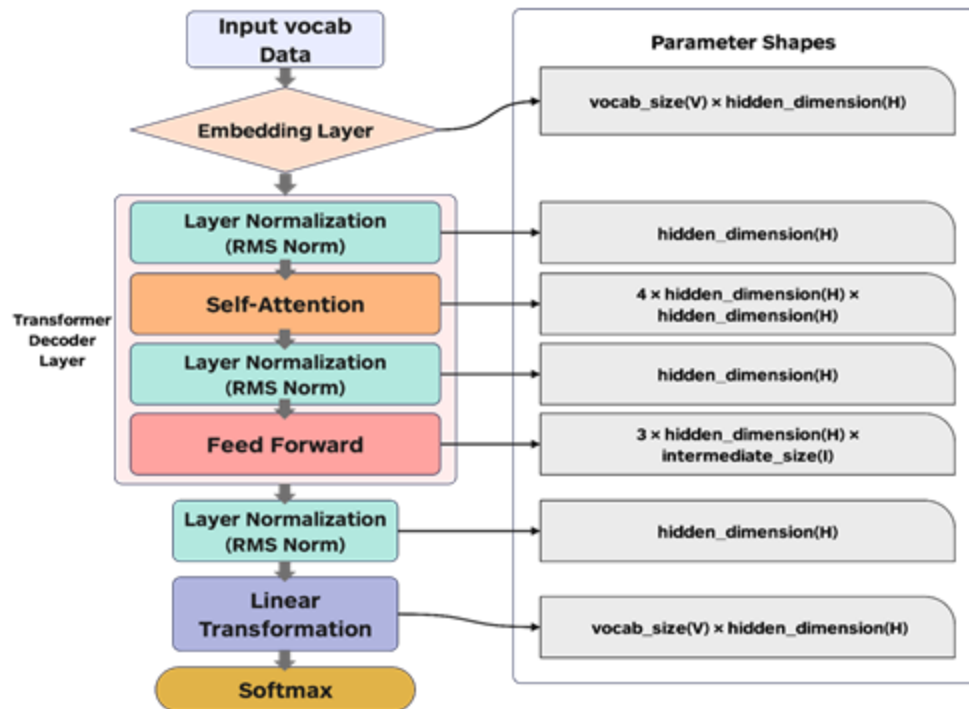
- SwiGLU helps the model capture more complex patterns by selectively gating information
- Swish is smoother than traditional activations ReLU



# Summary



# Scaling Up: Where is the Potential Bottleneck?

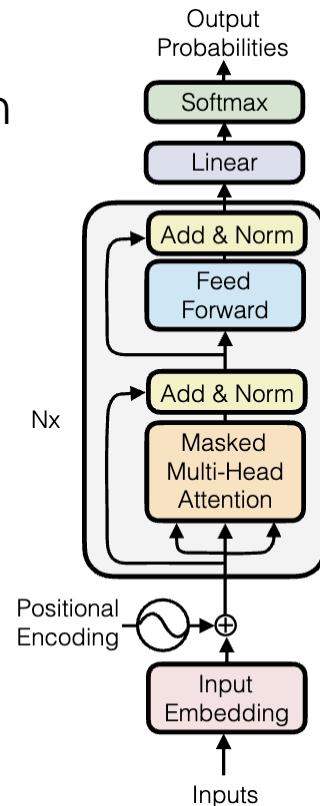


In PA3, you will implement this function 😊

# Connecting the Dots: Compute/Comm characteristic of LLMs

Key characteristics: compute, memory, communication

- calculate the number of parameters of an LLM?
- calculate the flops needed to train an LLM?
- calculate the memory needed to train an LLM?

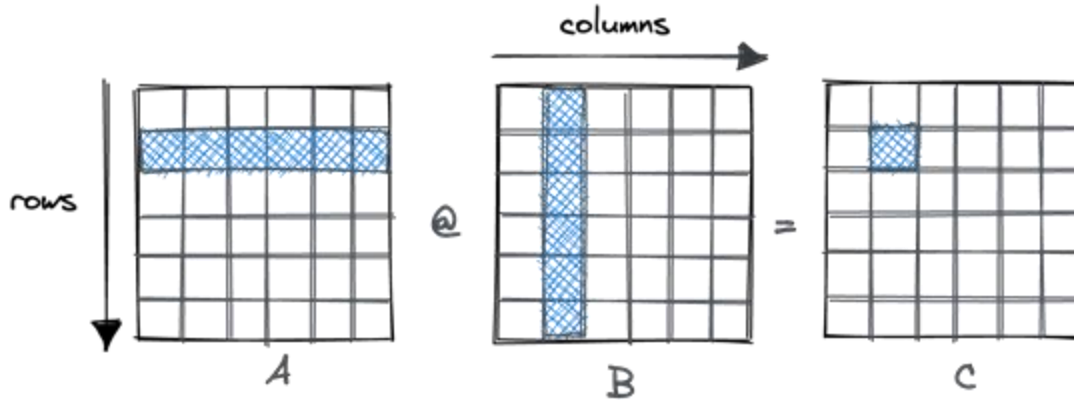


# Estimate the Compute: FLOPs

The FLOPs for multiplying two matrices of dimensions  $m \times n$  and  $n \times h$  can be calculated as follows:

$$\text{FLOPs} = m \times h \times (2n - 1)$$

So the total number of FLOPs is roughly  $\text{FLOPs} \approx 2m \times n \times h$



# LLama 2 7B Flops Forward Calculation (Training)

Hyperparameters:

Batch size:  $b$

Sequence length:  $s$

The number of attention heads:  $n$

Hidden state size of one head:  $d$

Hidden state size:  $h$  ( $h = n * d$ )

SwiGLU proj dim:  $i$

Vocab size:  $v$



Input:

$X$

Output Shape:

$(b, s, h)$

FLOPs

0

Self Attention:

$XW_Q, XW_K, XW_V$

$(b, s, h)$

$3 * 2bsh^2$

RoPE

$(b, n, s, d)$

$3bsnd$

$P = \text{Softmax}(QK^T/\sqrt{d})$

$(b, n, s, s)$

$2bs^2nd + 3bs^2n$

$PV$

$(b, n, s, d)$

$2bs^2nd$

$AW_O$

$(b, s, h)$

$2bsh^2$

Residual Connection:

$(b, s, h)$

$bsh$

Batch size:  $b$

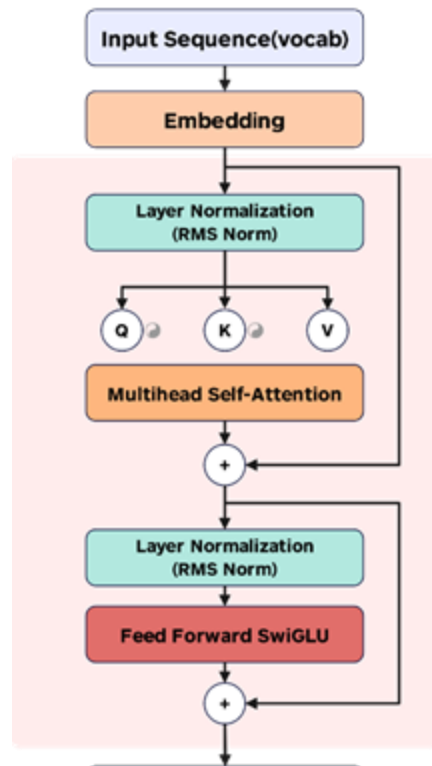
Sequence length:  $s$

# of attention heads:  $n$

Hidden state dim of one

head:  $d$

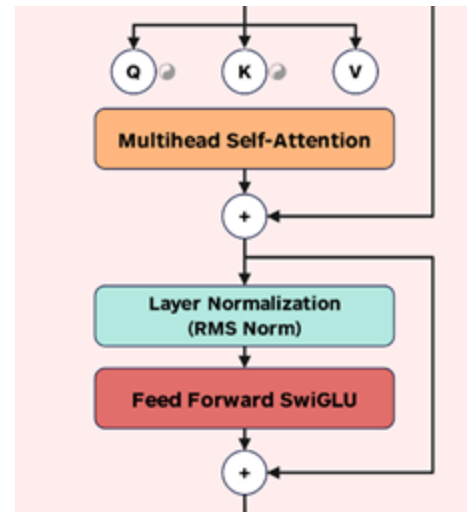
Hidden state dim:  $h$



Output from Self Attn:	Output Shape:	FLOPs
X	(b, s, h)	0
Feed-Forward SwiGLU:		
$XW_{\text{gate}}, XW_{\text{up}}$	(b, s, i)	$2 * 2bshi$
Swish Activation	(b, s, i)	$4bsi$
Element-wise *	(b, s, i)	$bsi$
$XW_{\text{down}}$	(b, s, h)	$2bshi$
RMS Norm:		
	(b, s, h)	$4bsh + 2bs$

$$\text{SwiGLU}(x) = \text{Swish}(xW_1 + b_1) \odot (xW_2 + b_2)$$

Batch size:  $b$   
 Sequence length:  $s$   
 Hidden state dim:  $h$   
 SwiGLU proj dim:  $i$



1. Calculate Root Mean Square:

$$\bullet \text{ RMS}(x) = \sqrt{\frac{1}{d} \sum_{i=1}^d x_i^2}$$

2. Normalize:

$$\bullet \text{ RMSNorm}(x) = \frac{x}{\text{RMS}(x) + \epsilon} \cdot \gamma$$

# LLama 2 7B Flops Forward (Training)

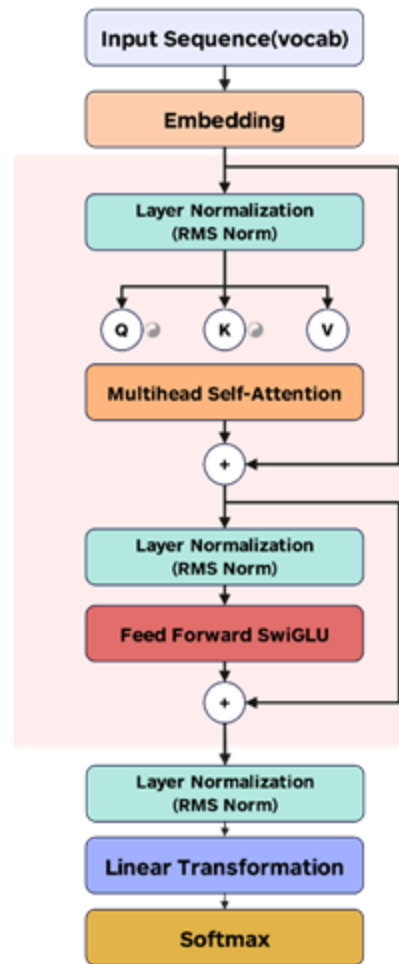
Total Flops  $\approx$  #num\_layers \* (Attention block + SwiGLU block)

+ Prediction head

= #num\_layers \* ( $6bsh^2 + 4bs^2h + 3bs^2n + 2bsh^2$ )

+ #num\_layers ( 6bshi)

+ 2 bshv



# LLama 2 7B Flops Forward Calculation (Training)

Hyperparameters:

Batch size:  $b=1$

Sequence length:  $s=4096$

The number of attention heads:  $n=32$

Hidden state size of one head:  $d=128$

Hidden state size:  $h=4096$

SwiGLU proj dim:  $i=11008$

Vocab size:  $v=32000$

The number of layers:  $N=32$

$$\text{Total Flops} \approx N * (6bsh^2 + 4bs^2h + 3bs^2n + 2bsh^2)$$

$$+ N (6bshi)$$

$$+ 2 bshv$$

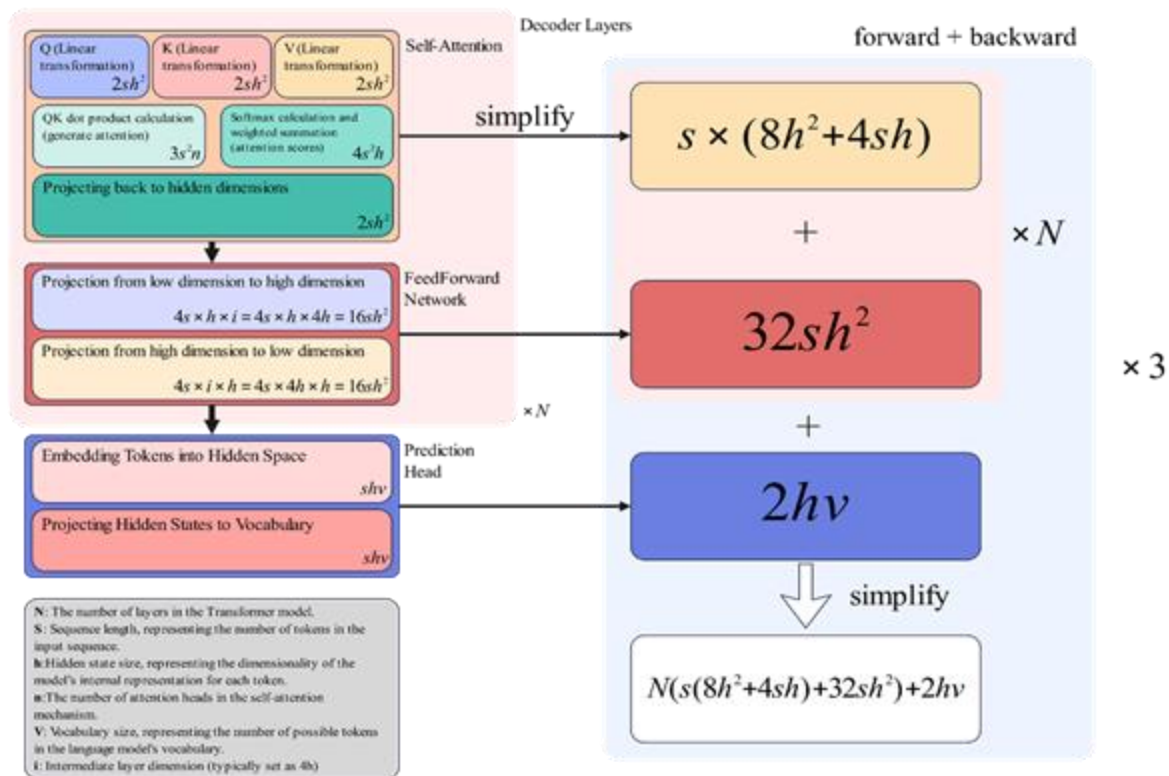
$$\approx 63 \text{ TFLOPs}$$

# Flops Distribution

## Part 1: Training Computational Costs

- **Total Training TeraFLOPs:** 192.17 TFLOPs
- **FLOP Distribution by Layer:**
  - **Embedding Layer:** 1.676%
  - **Normalization:** 0.007%
  - **Residual:** 0.003%
  - **Attention:** 41.276%
  - **MLP (Multi-Layer Perceptron):** 55.361%
  - **Linear:** 1.676%

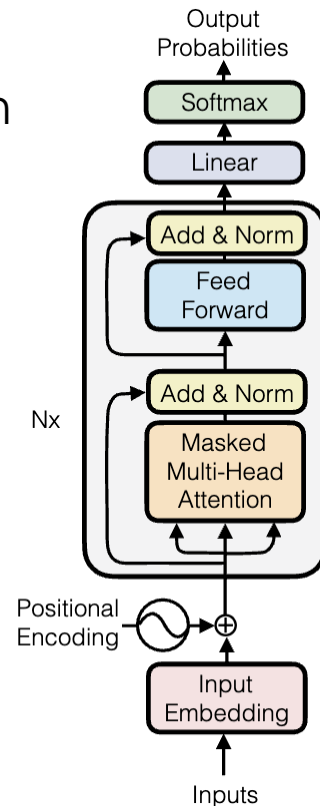
# Scaling Up: Where is the Potential Bottleneck?



# Connecting the Dots: Compute/Comm characteristic of LLMs

Key characteristics: compute, memory, communication

- calculate the number of parameters of an LLM?
- calculate the flops needed to train an LLM?
- calculate the memory needed to train an LLM?



### Composition of Memory Usage (Training)

Model Weights

Intermediate Action Value

Optimizer States

Weight Gradients + Activation Gradients



## Llama2-7b Mix Precision(16bit-32bit)

**b:** Batch size

**s:** Max sequence Length

**h:** Hidden Dimension

**i:** Intermediate Size

**n:** Number of heads

**d:** Head Dimension ( $n \times d = h$ )

**v:** Vocabulary Size

## Llama2-7b Mix Precision(16bit-32bit)

X input

(b, s, v)

**b:** Batch size

**s:** Max sequence Length

**h:** Hidden Dimension

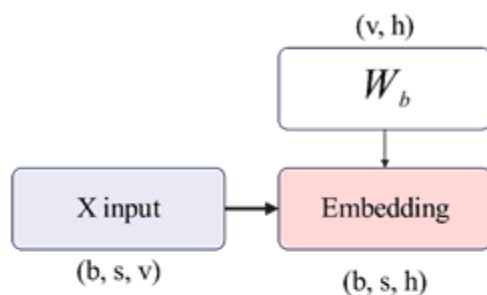
**i:** Intermediate Size

**n:** Number of heads

**d:** Head Dimension ( $n \times d = h$ )

**v:** Vocabulary Size

## Llama2-7b Mix Precision(16bit-32bit)



**b:** Batch size

**s:** Max sequence Length

**h:** Hidden Dimension

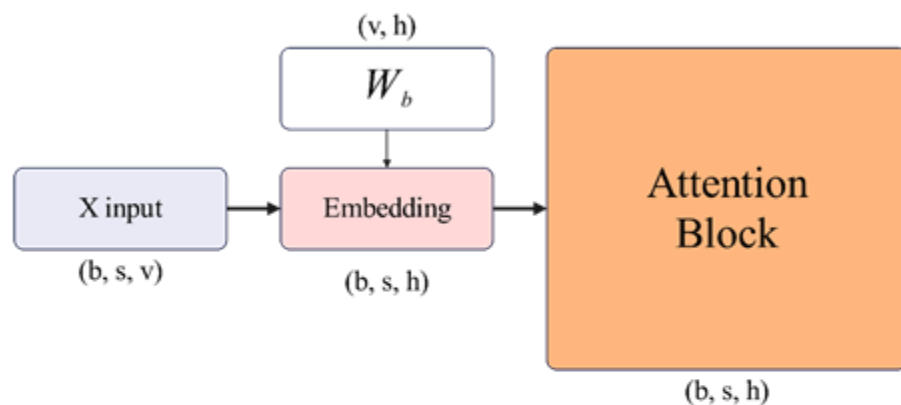
**i:** Intermediate Size

**n:** Number of heads

**d:** Head Dimension ( $n \times d = h$ )

**v:** Vocabulary Size

## Llama2-7b Mix Precision(16bit-32bit)



**b:** Batch size

**s:** Max sequence Length

**h:** Hidden Dimension

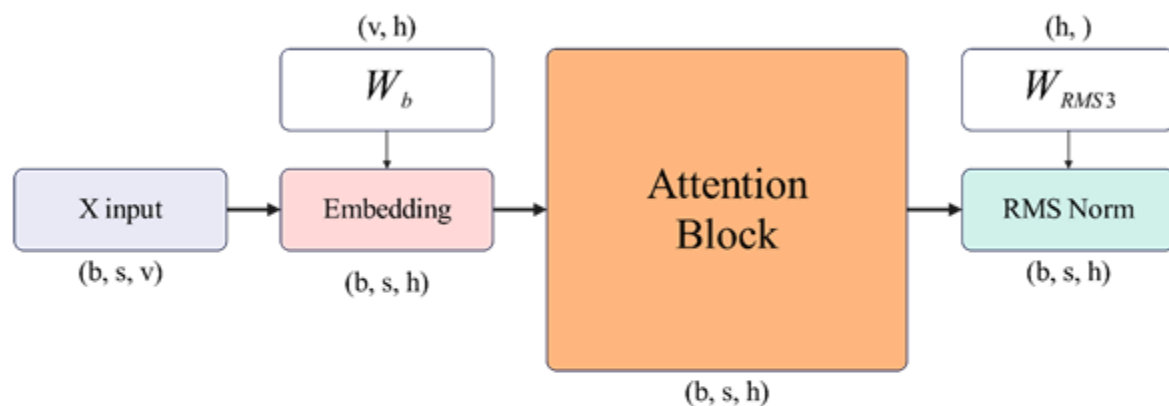
**i:** Intermediate Size

**n:** Number of heads

**d:** Head Dimension ( $n \times d = h$ )

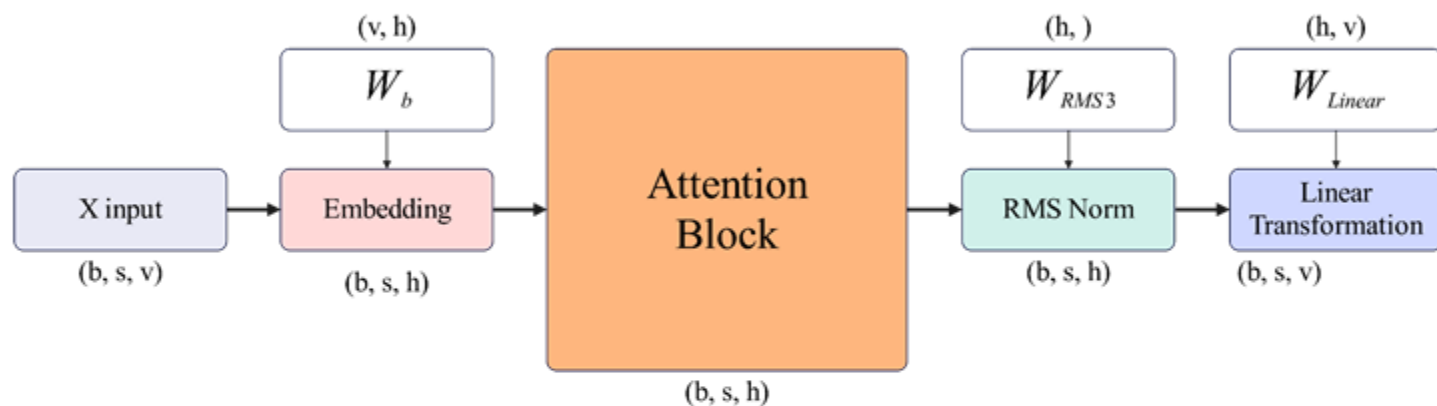
**v:** Vocabulary Size

## Llama2-7b Mix Precision(16bit-32bit)



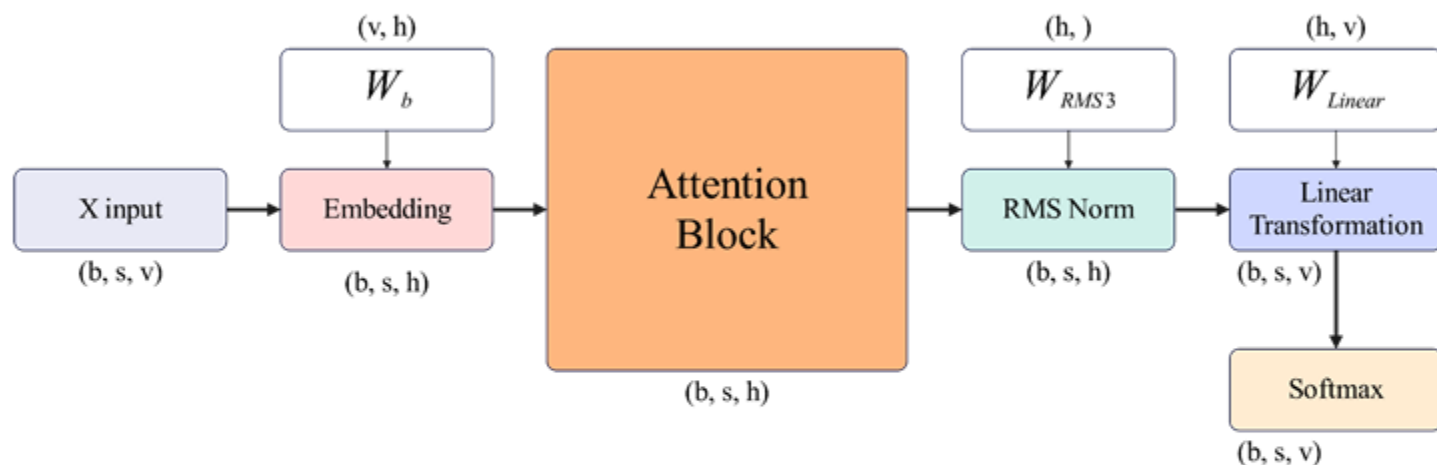
**b:** Batch size  
**s:** Max sequence Length  
**h:** Hidden Dimension  
**i:** Intermediate Size  
**n:** Number of heads  
**d:** Head Dimension ( $n \times d = h$ )  
**v:** Vocabulary Size

## Llama2-7b Mix Precision(16bit-32bit)



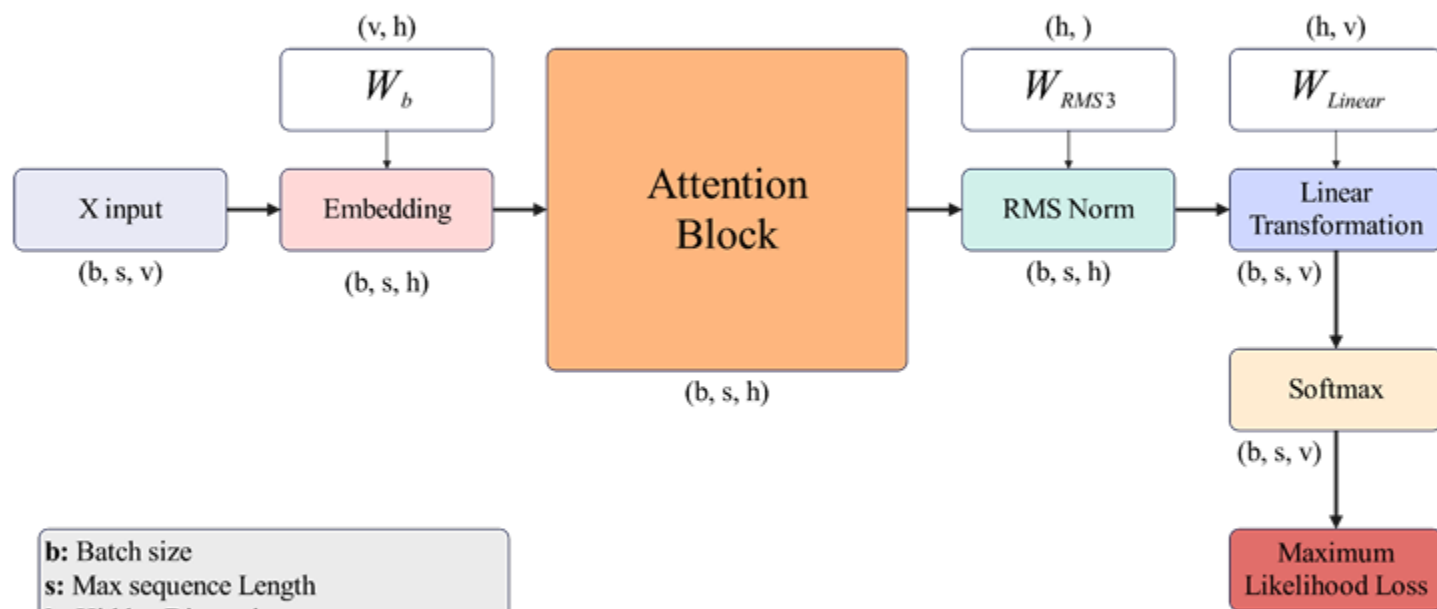
**b:** Batch size  
**s:** Max sequence Length  
**h:** Hidden Dimension  
**i:** Intermediate Size  
**n:** Number of heads  
**d:** Head Dimension ( $n \times d = h$ )  
**v:** Vocabulary Size

## Llama2-7b Mix Precision(16bit-32bit)



**b:** Batch size  
**s:** Max sequence Length  
**h:** Hidden Dimension  
**i:** Intermediate Size  
**n:** Number of heads  
**d:** Head Dimension ( $n \times d = h$ )  
**v:** Vocabulary Size

## Llama2-7b Mix Precision(16bit-32bit)



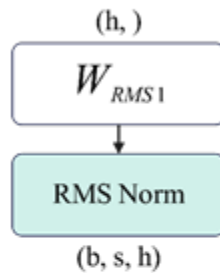
**b:** Batch size  
**s:** Max sequence Length  
**h:** Hidden Dimension  
**i:** Intermediate Size  
**n:** Number of heads  
**d:** Head Dimension ( $n \times d = h$ )  
**v:** Vocabulary Size



## Llama2-7b Attention Block (Self-Attention)

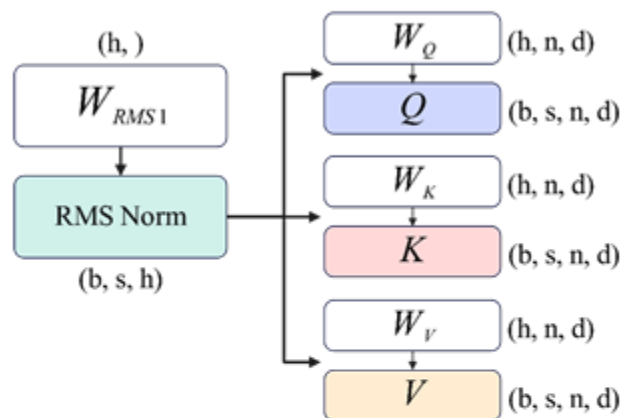
**b:** Batch size  
**s:** Max sequence Length  
**h:** Hidden Dimension  
**i:** Intermediate Size  
**n:** Number of heads  
**d:** Head Dimension ( $n \times d = h$ )  
**v:** Vocabulary Size

## Llama2-7b Attention Block (Self-Attention)



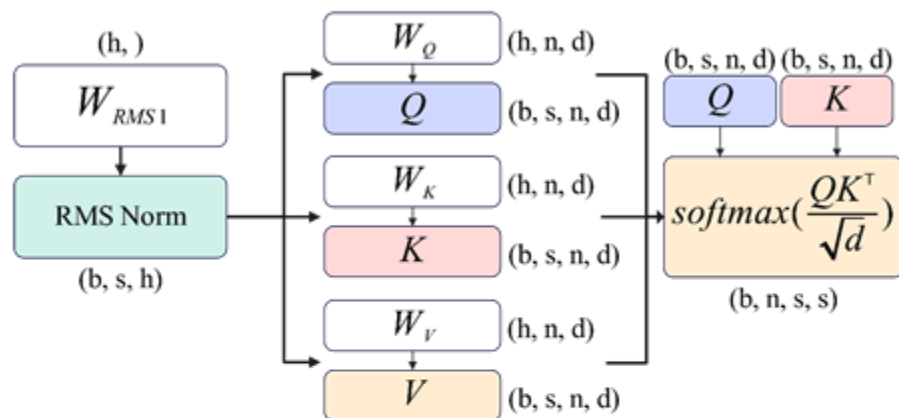
**b:** Batch size  
**s:** Max sequence Length  
**h:** Hidden Dimension  
**i:** Intermediate Size  
**n:** Number of heads  
**d:** Head Dimension ( $n \times d = h$ )  
**v:** Vocabulary Size

## Llama2-7b Attention Block (Self-Attention)



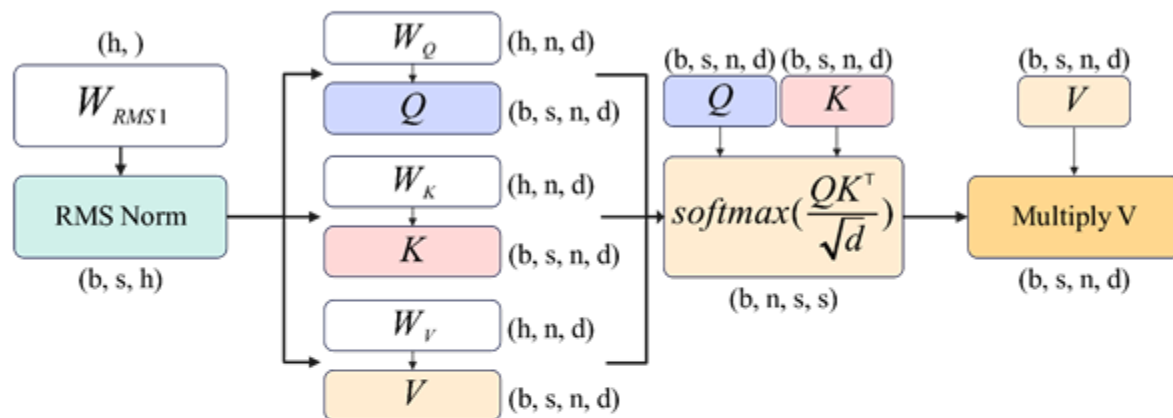
**b:** Batch size  
**s:** Max sequence Length  
**h:** Hidden Dimension  
**i:** Intermediate Size  
**n:** Number of heads  
**d:** Head Dimension ( $n \times d = h$ )  
**v:** Vocabulary Size

## Llama2-7b Attention Block (Self-Attention)



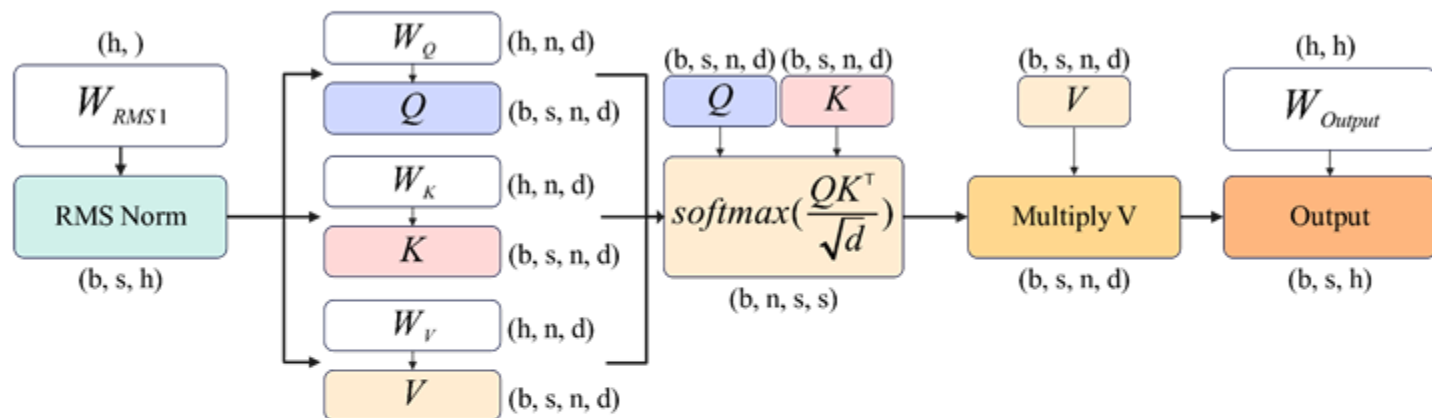
**b:** Batch size  
**s:** Max sequence Length  
**h:** Hidden Dimension  
**i:** Intermediate Size  
**n:** Number of heads  
**d:** Head Dimension ( $n \times d = h$ )  
**v:** Vocabulary Size

## Llama2-7b Attention Block (Self-Attention)



**b:** Batch size  
**s:** Max sequence Length  
**h:** Hidden Dimension  
**i:** Intermediate Size  
**n:** Number of heads  
**d:** Head Dimension ( $n \times d = h$ )  
**v:** Vocabulary Size

## Llama2-7b Attention Block (Self-Attention)

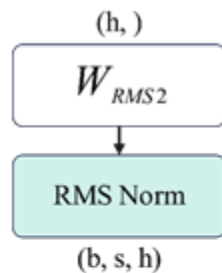


**b:** Batch size  
**s:** Max sequence Length  
**h:** Hidden Dimension  
**i:** Intermediate Size  
**n:** Number of heads  
**d:** Head Dimension ( $n \times d = h$ )  
**v:** Vocabulary Size

## Llama2-7b Attention Block (FeedForward)

**b:** Batch size  
**s:** Max sequence Length  
**h:** Hidden Dimension  
**i:** Intermediate Size  
**n:** Number of heads  
**d:** Head Dimension ( $n \times d = h$ )  
**v:** Vocabulary Size

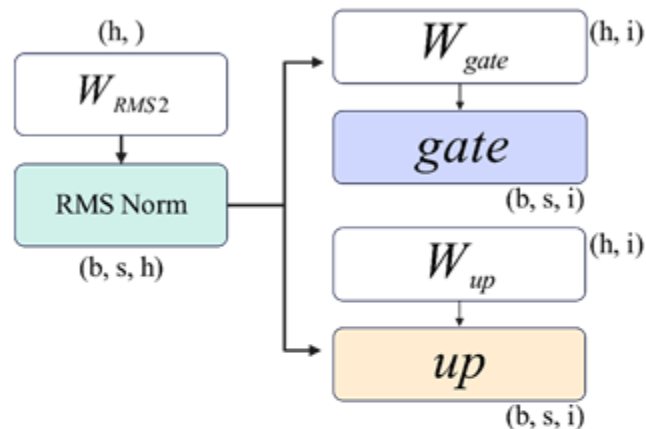
## Llama2-7b Attention Block (FeedForward)



**b:** Batch size  
**s:** Max sequence Length  
**h:** Hidden Dimension  
**i:** Intermediate Size  
**n:** Number of heads  
**d:** Head Dimension ( $n \times d = h$ )  
**v:** Vocabulary Size

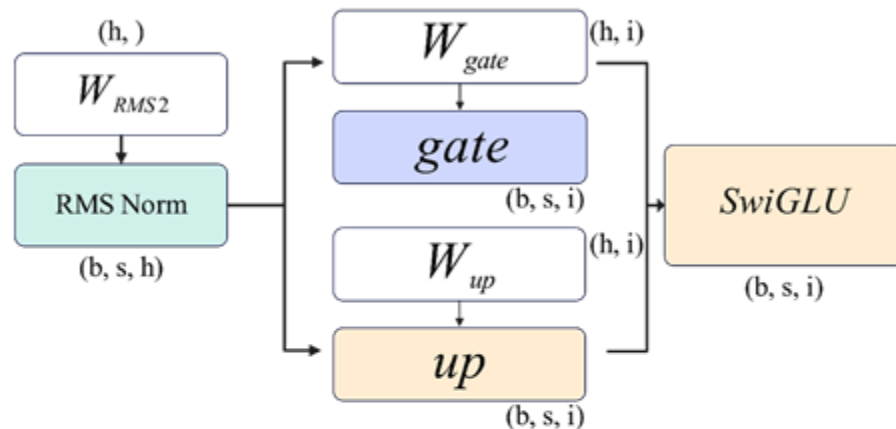


## Llama2-7b Attention Block (FeedForward)



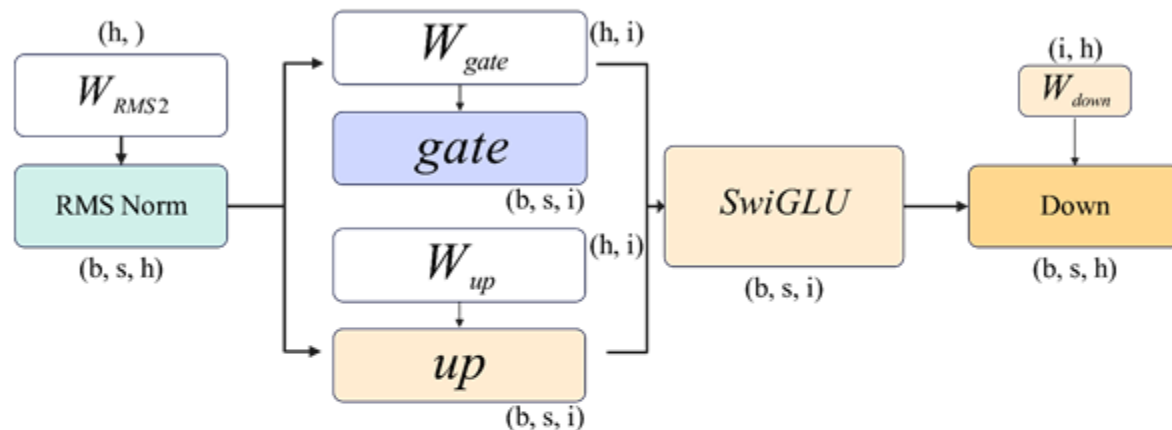
**b:** Batch size  
**s:** Max sequence Length  
**h:** Hidden Dimension  
**i:** Intermediate Size  
**n:** Number of heads  
**d:** Head Dimension ( $n \times d = h$ )  
**v:** Vocabulary Size

## Llama2-7b Attention Block (FeedForward)



**b:** Batch size  
**s:** Max sequence Length  
**h:** Hidden Dimension  
**i:** Intermediate Size  
**n:** Number of heads  
**d:** Head Dimension ( $n \times d = h$ )  
**v:** Vocabulary Size

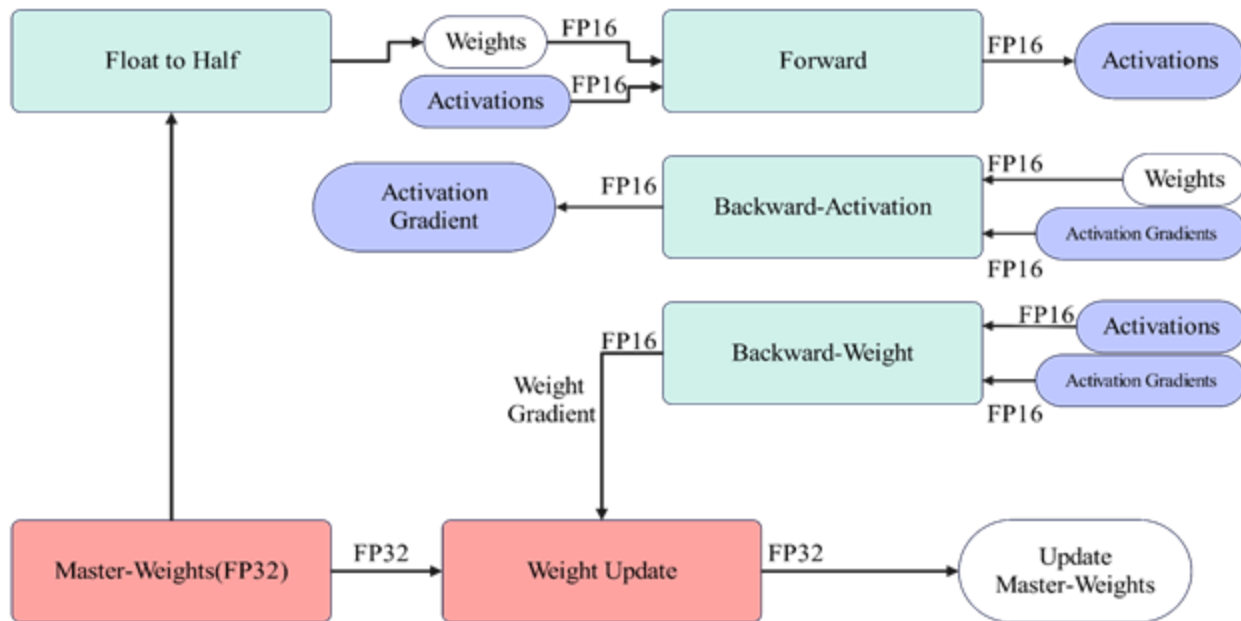
## Llama2-7b Attention Block (FeedForward)



**b:** Batch size  
**s:** Max sequence Length  
**h:** Hidden Dimension  
**i:** Intermediate Size  
**n:** Number of heads  
**d:** Head Dimension ( $n \times d = h$ )  
**v:** Vocabulary Size

# Optimizer States: 16M

## Mixed-Precision



# Large Language Models

- Transformers, Attentions
- **Scaling Law**
  - MoE
- Connecting the dots: Training Optimizations
  - Flash attention
  - Long context, parallelism
- Serving and inference optimization
  - Continuous batching and Paged attention
  - Speculative decoding (Guest Lecture)
- Connecting the dots: Deepseek-v3
- Hot topics

## Some Observations

- compute is a function of:  $h, i, b$
- #parameter is a function of:  $h, i$
- Hence: compute correlates with #parameters
  - more parameters, more compute
  - more data, more compute (of course)
- Problem: we have limited compute (\$)
- how should we allocate our limited resources:
  - Train models longer vs train bigger models?
  - Collect more data vs get more GPUs?

# Motivation of Scaling Laws

- We want to know:
  - how large a model should we train...
  - How many data should we use...
  - To achieve a given performance...
  - Subject to a compute budget (\$)?

How do we do that in traditional ML: data scaling law

**Input:**  $x_1 \dots x_n \sim N(\mu, \sigma^2)$

**Task:** estimate the average as  $\hat{\mu} = \frac{\sum_i x_i}{n}$

**What's the error?** By standard arguments..

$$E[(\hat{\mu} - \mu)^2] = \frac{\sigma^2}{n}$$

**This is a scaling law!!**

$$\log(\text{Error}) = -\log n + 2 \log \sigma$$

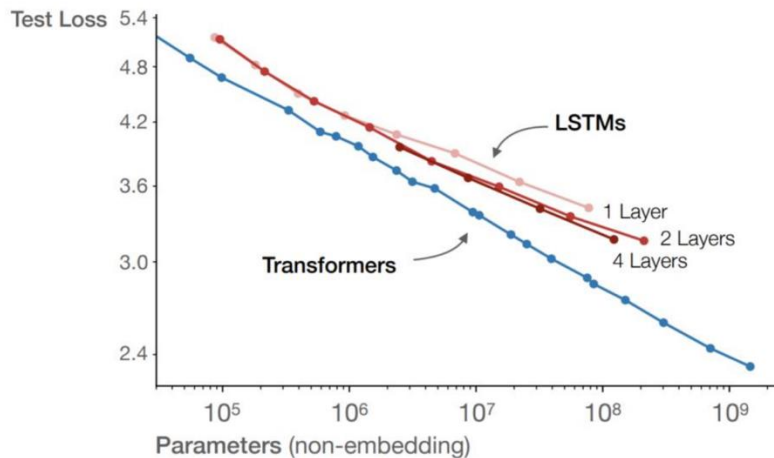
- Can we do this for transformers LLMs?

More generally, any polynomial rate  $1/n^\alpha$  is a scaling law



# Transformers vs LSTMs

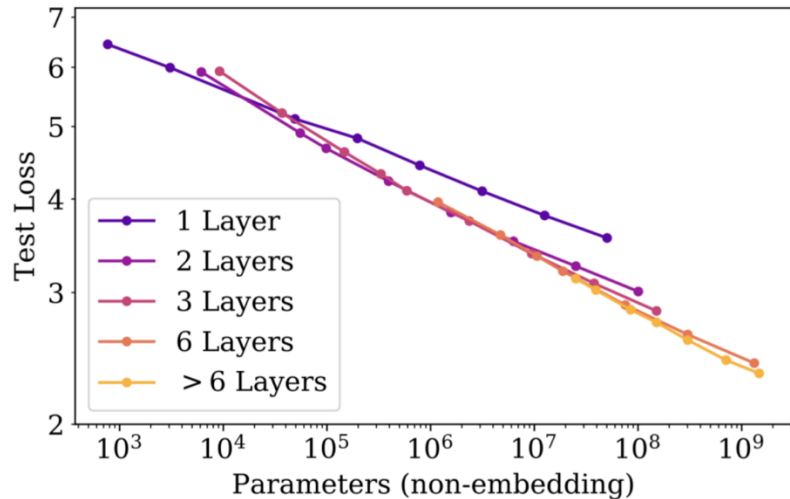
- Q: Are transformers better than LSTMs?
  - Brute force way: spend tens of millions to train a LSTM GPT-3
- Scaling law way:



[Kaplan+ 2021]

# Number of Layers

- Does depth or width make a huge difference?
  - 1 vs 2 layers makes a huge difference.
  - More layers have diminishing returns below  $10^7$  params



# The Scaling law way

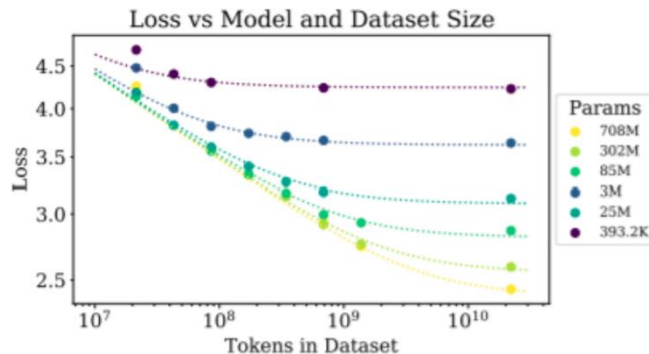
- Approach:
  - Train a few smaller models
  - Establish a scaling law (LSTM vs. transformers)
  - Select optimal hyperparam based on the scaling law prediction.
- Rationale
  - The effect of hyperparameters on big LMs can be predicted before training!
    - Optimizer choice
    - Model Depth
    - Architecture choice

## Back to our problem:

- how large a model should we train...
  - How many data should we use...
  - To achieve a given performance...
  - Subject to a compute budget?
- 
- Approach: estimate a law between model size data joint scaling

# Model size data joint scaling

- Do we need more data or bigger model
  - Clearly, lots of data is wasted on small models
- Joint data-model scaling laws describe how the two relate



From Rosenfeld+ 2020,

$$Error = n^{-\alpha} + m^{-\beta} + C$$

From Kaplan+ 2021

$$Error = [m^{-\alpha} + n^{-1}]^{\beta}$$

Provides surprisingly good fits to model-data joint error.

# Compute Trade-offs

- Q: what about other resources? Compute vs. performance?
- For a fixed compute budget...
  - Big models that's undertrained vs small model that's well trained?
  - Solving the following optimization?

$$N_{opt}(C), D_{opt}(C) = \underset{N, D \text{ s.t. } \text{FLOPs}(N, D) = C}{\operatorname{argmin}} L(N, D).$$

# Approach: empirical scaling law

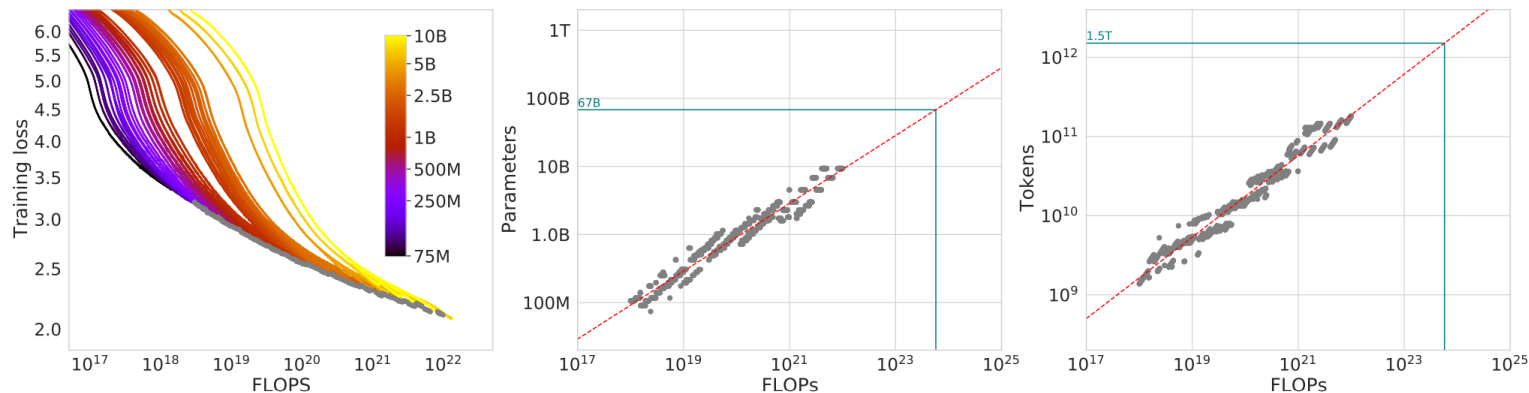


Figure 2 | **Training curve envelope.** On the **left** we show all of our different runs. We launched a range of model sizes going from 70M to 10B, each for four different cosine cycle lengths. From these curves, we extracted the envelope of minimal loss per FLOP, and we used these points to estimate the optimal model size (**center**) for a given compute budget and the optimal number of training tokens (**right**). In green, we show projections of optimal model size and training token count based on the number of FLOPs used to train *Gopher* ( $5.76 \times 10^{23}$ ).

Today's SoTA Law

$$L(N, D) = \frac{406.4}{N^{0.34}} + \frac{410.7}{D^{0.29}} + 1.69$$



# Summary

- Scaling law: the physics behind LLMs
- Scaling law also represents a research approach transition:
  - Rigorous theoretical analysis -> empirical laws
  - Exploration of different model architectures -> Scaling transformers
- Due to scaling law: ML systems become essential

## PA3: Q3

You already know:

- How to estimate the number of parameters of an LLM?
  - How to estimate the flops needed to train an LLM?
  - How to estimate the memory needed to train a transformer?
- 
- We will give you a scaling law and compute budget
    - Task: design your optimal LLM

## Next Lecture: What is MoE

- Superficially: experts
- Essentially: a model with a better scaling law.