



<https://hao-ai-lab.github.io/dsc204a-f25/>

# DSC 204A: Scalable Data Systems

## Fall 2025

---

Staff

Instructor: Hao Zhang

TAs: Mingjia Huo, Yuxuan Zhang

 [@haozhangml](https://twitter.com/haozhangml)

 [@haoailab](https://twitter.com/haoailab)

 [haozhang@ucsd.edu](mailto:haozhang@ucsd.edu)

# Logistics

- All 4 Guest Lectures Complete
  - Eugune Wu: classic DB + HCI researcher
  - Shreya Shankar: Modern DB + HCI researcher
- Andrey Cheng:
  - DB researcher, but now working on LLM for optimizing DB
- Junchen Jiang: networking, but want to propose a new type of DB
- **Fall 2025 Student Evaluations of Teaching were sent**
  - **Again: if 80% of you finish the evaluation, all will get 2 bonus points.**

# Logistics

## **Fall 2025 Student Evaluations of Teaching were sent**

- Again: if 80% of you finish the evaluation, all will get 2 bonus points.

We might need 1 – 2 extra lectures (beyond scheduled) to compensate holiday interruptions

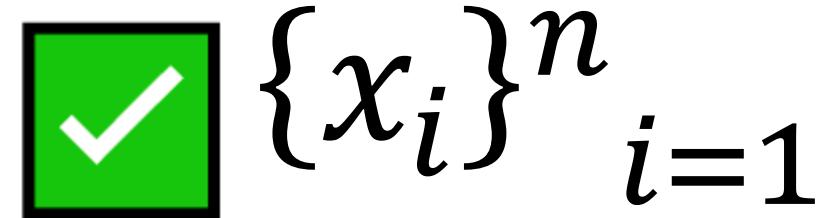
- Will decide depending on our progress
- If happen – will on Zoom

Exam:

- all MCQ
- TA will hold a recitation before exam
- Date and time TBD

# High-level Picture

Data



Model

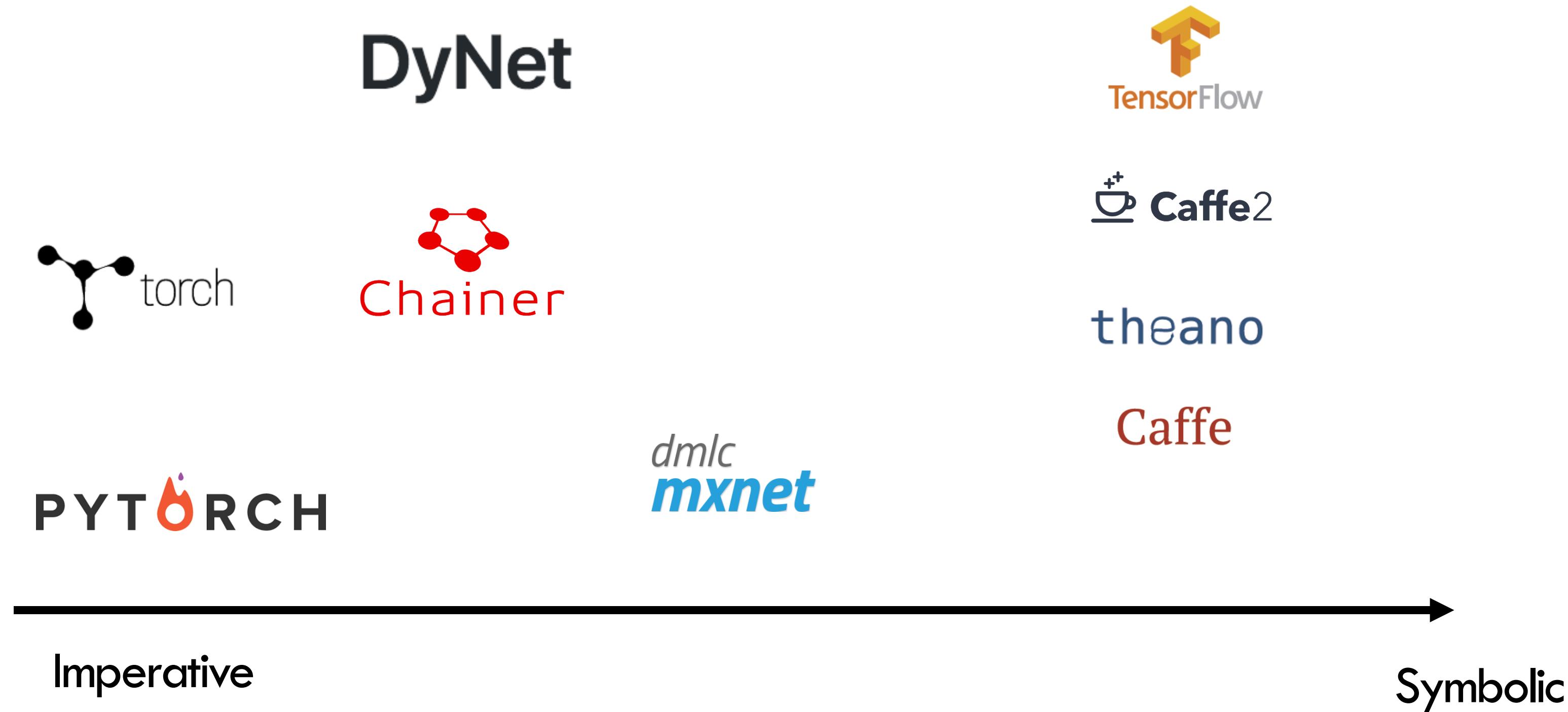
Math primitives  
(mostly matmul)

? A repr that expresses the computation using primitives

Compute

? Make them run on (clusters of ) different kinds of hardware

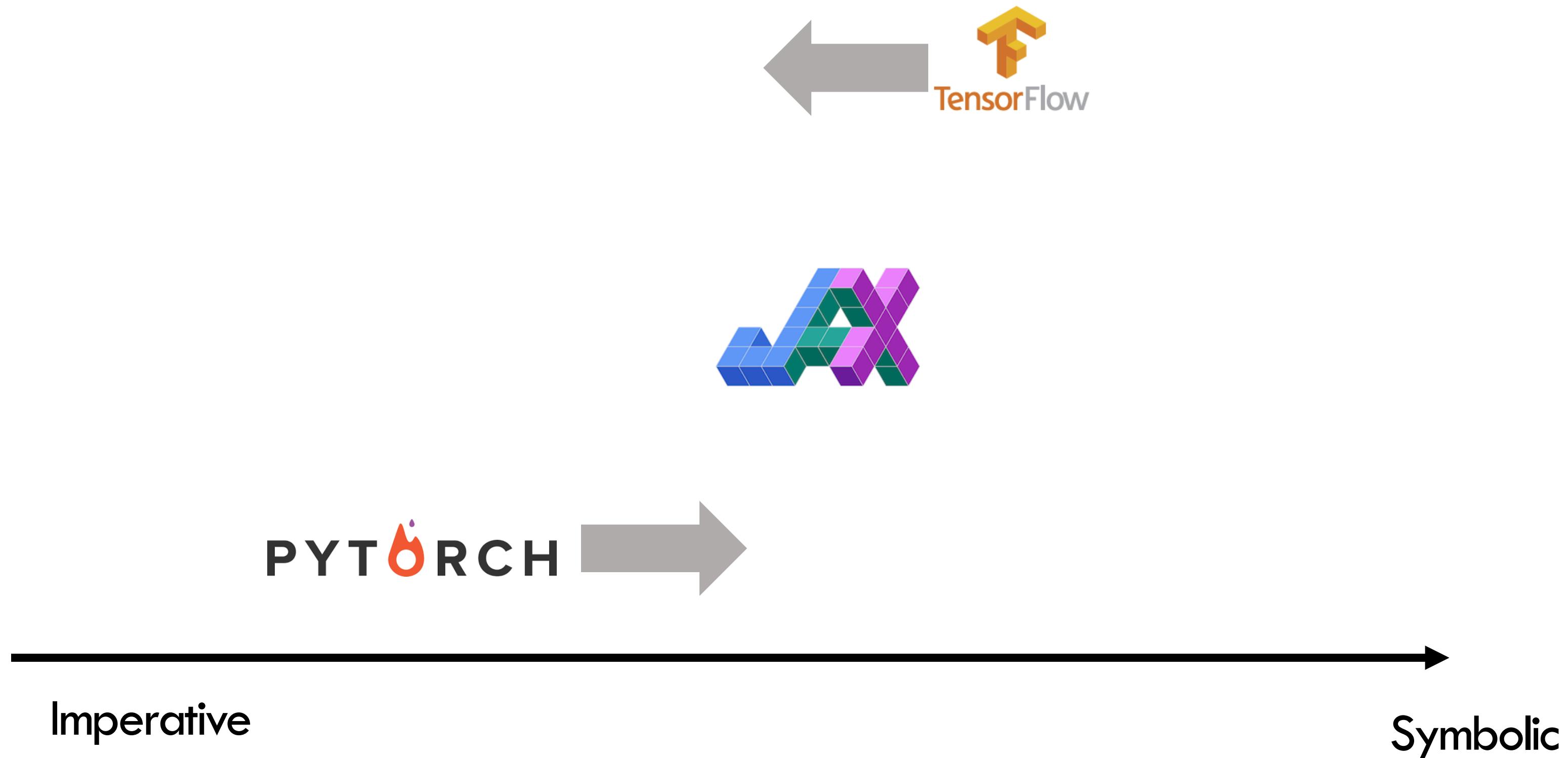
# Symbolic vs. Imperative (2016)



# Symbolic vs. Imperative (2024)



# Symbolic vs. Imperative (2024)



# Just-in-time (JIT) Compilation

- Ideally, we want define-and-run during \_\_\_\_\_
- We want define-then-run during \_\_\_\_\_
- Q: how can combine the best of both worlds?

```
x = torch.Tensor([3])
y = torch.Tensor([2])
z = x - y
loss = square(z)
loss.backward()
print(x.grad)
```

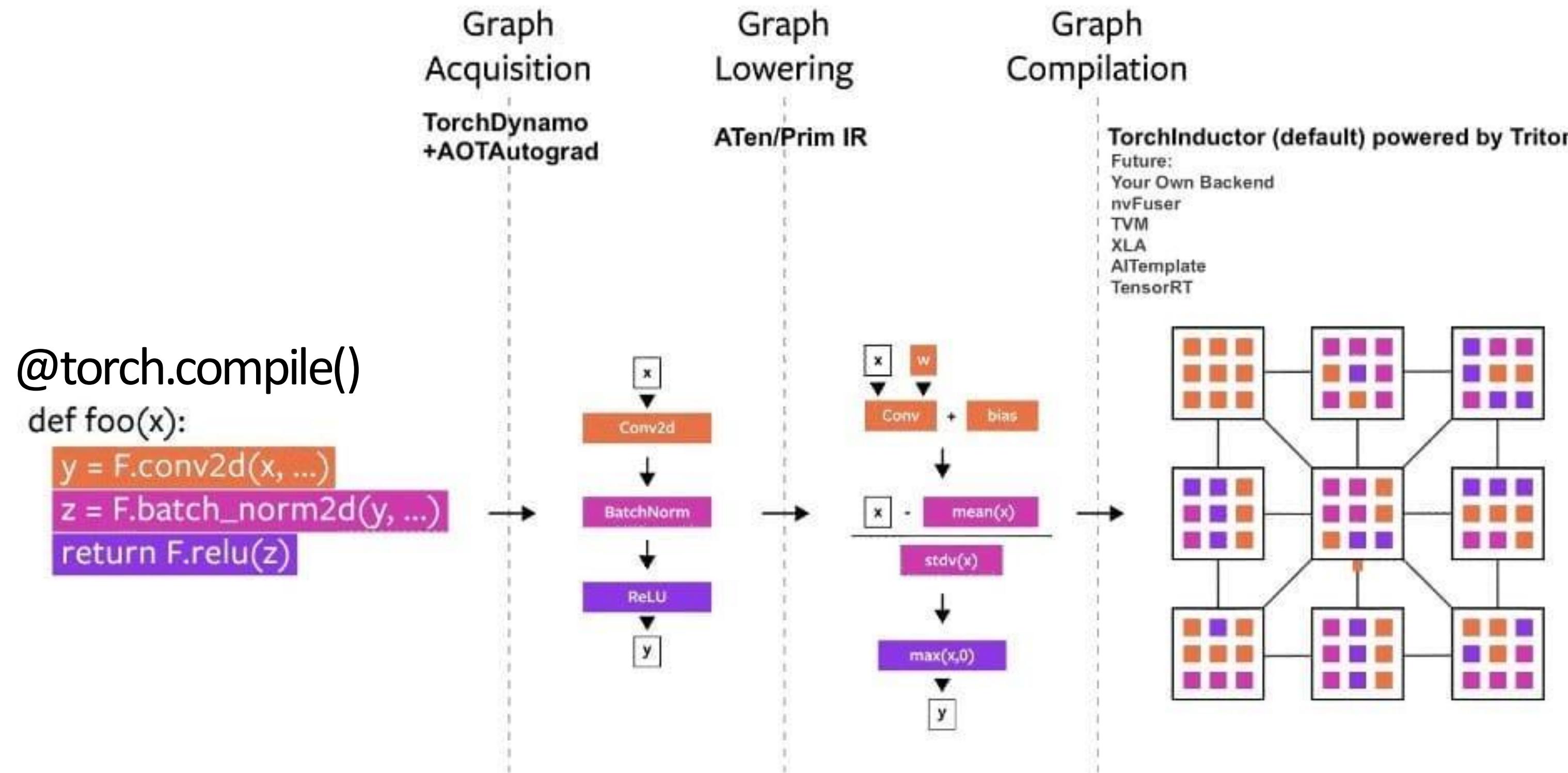
**Dev mode**

```
@torch.compile()

x = torch.Tensor([3])
y = torch.Tensor([2])
z = x - y
loss = square(z)
loss.backward()
print(x.grad)
```

**Deploy mode:**  
**Decorate `torch.compile()`**

# What happens behind the scene

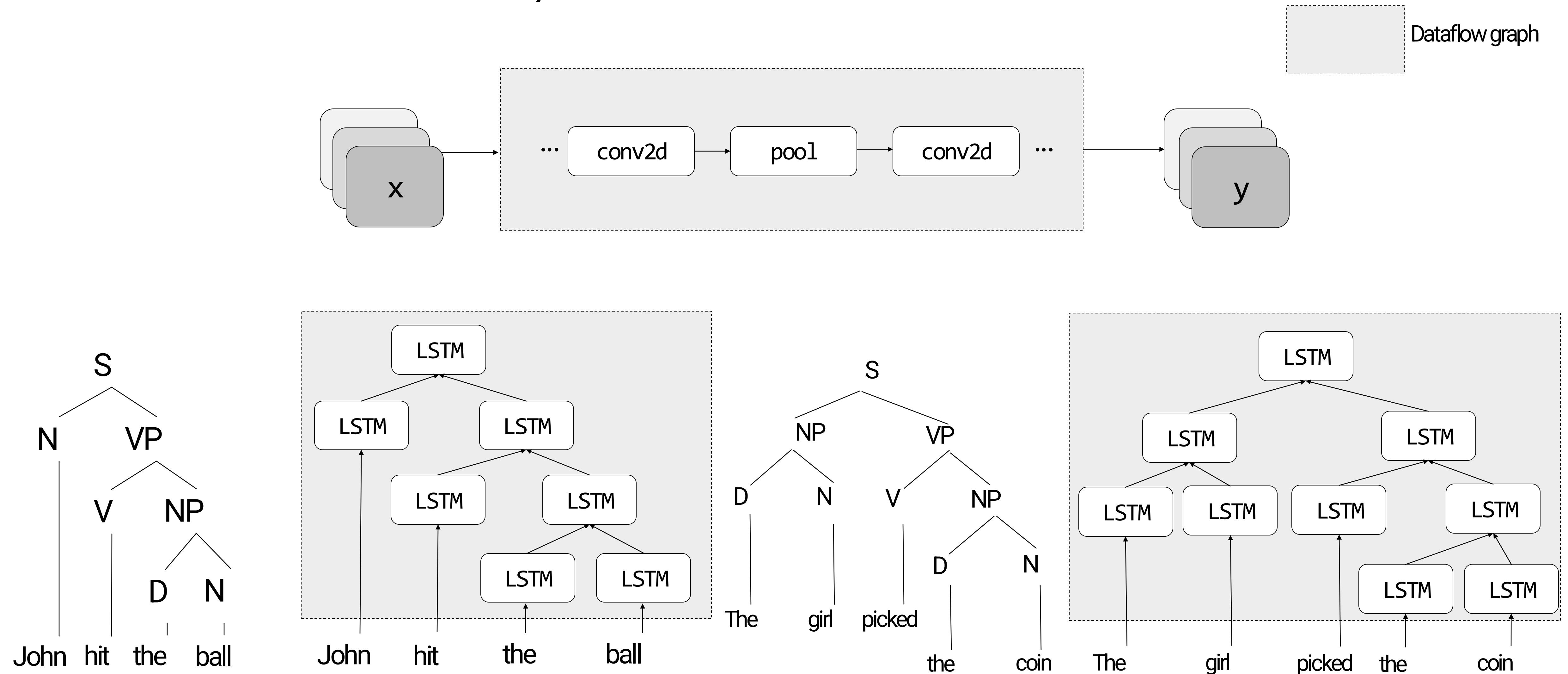


What is the problem of JIT?  
Requirements for static graphs

Q: What is the problem of JIT?

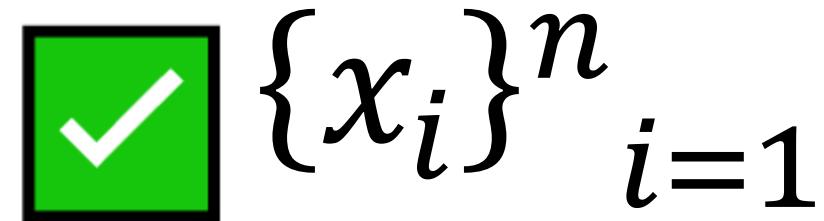
A: Requirements for static graphs

# Static Models vs. Dynamic Models



# High-level Picture

Data



Model

Math primitives  
(mostly matmul)

? A repr that expresses the computation using primitives

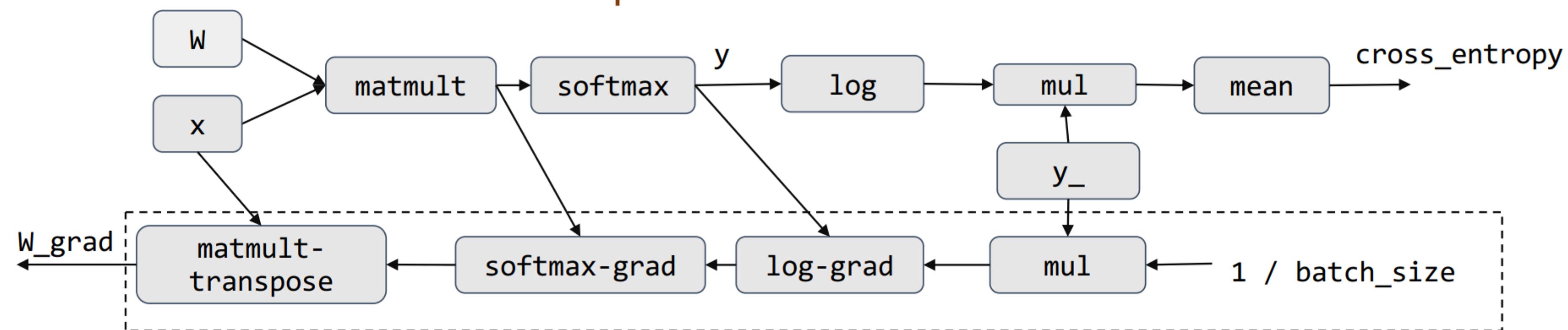
Compute

? Make them run on (clusters of ) different kinds of hardware

# What happens behind the Scene (Cond.)

```
w_grad = tf.gradients(cross_entropy, [W])[0]
```

Automatic Differentiation, more details in follow up lectures



# Expand it a Bit

A repr that expresses the computation using primitives

 A repr that expresses the **forward** computation using primitives

 A repr that expresses the **backward** computation using primitives

Recap: how to take derivative?

Given  $f(\theta)$ , what is  $\frac{\partial f}{\partial \theta}$  ?

## Instead, Symbolic Differentiation

Write down the formula, derive the gradient following PD rules

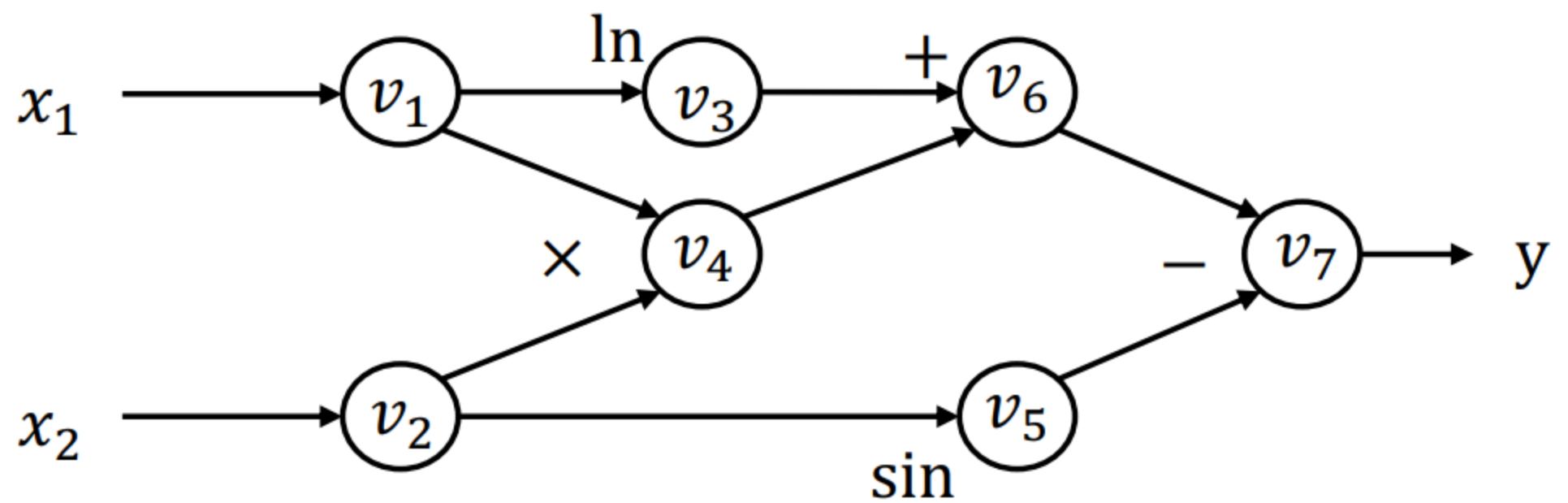
$$\frac{\partial(f(\theta) + g(\theta))}{\partial\theta} = \frac{\partial f(\theta)}{\partial\theta} + \frac{\partial g(\theta)}{\partial\theta}$$

$$\frac{\partial(f(\theta)g(\theta))}{\partial\theta} = g(\theta) \frac{\partial f(\theta)}{\partial\theta} + f(\theta) \frac{\partial g(\theta)}{\partial\theta}$$

$$\frac{\partial(f(g(\theta)))}{\partial\theta} = \frac{\partial f(g(\theta))}{\partial g(\theta)} \frac{\partial g(\theta)}{\partial\theta}$$

# Map autodiff rules to computational graph

$$y = f(x_1, x_2) = \ln(x_1) + x_1 x_2 - \sin x_2$$



Forward evaluation trace

$$v_1 = x_1 = 2$$

$$v_2 = x_2 = 5$$

$$v_3 = \ln v_1 = \ln 2 = 0.693$$

$$v_4 = v_1 \times v_2 = 10$$

$$v_5 = \sin v_2 = \sin 5 = -0.959$$

$$v_6 = v_3 + v_4 = 10.693$$

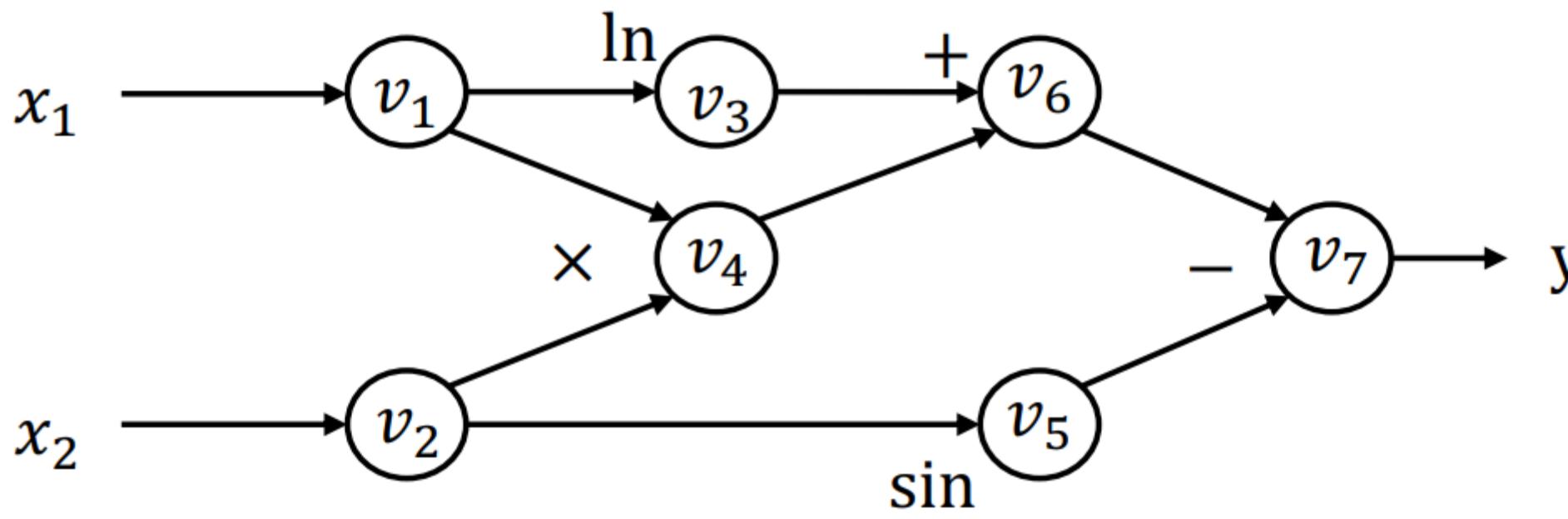
$$v_7 = v_6 - v_5 = 10.693 + 0.959 = 11.652$$

$$y = v_7 = 11.652$$

- Q: Calculate the value of  $\frac{\partial y}{\partial x_1}$
- A: use PD and chain rules

# Reverse Mode AD

$$y = f(x_1, x_2) = \ln(x_1) + x_1 x_2 - \sin x_2$$



Forward evaluation trace

$$v_1 = x_1 = 2$$

$$v_2 = x_2 = 5$$

$$v_3 = \ln v_1 = \ln 2 = 0.693$$

$$v_4 = v_1 \times v_2 = 10$$

$$v_5 = \sin v_2 = \sin 5 = -0.959$$

$$v_6 = v_3 + v_4 = 10.693$$

$$v_7 = v_6 - v_5 = 10.693 + 0.959 = 11.652$$

$$y = v_7 = 11.652$$

- Define adjoint  $\bar{v}_i = \frac{\partial y}{\partial v_i}$
- We then compute each  $\bar{v}_i$  in the reverse topological order of the graph

$$\bar{v}_7 = \frac{\partial y}{\partial v_7} = 1$$

$$\bar{v}_6 = \bar{v}_7 \frac{\partial v_7}{\partial v_6} = \bar{v}_7 \times 1 = 1$$

$$\bar{v}_5 = \bar{v}_7 \frac{\partial v_7}{\partial v_5} = \bar{v}_7 \times (-1) = -1$$

$$\bar{v}_4 = \bar{v}_6 \frac{\partial v_6}{\partial v_4} = \bar{v}_6 \times 1 = 1$$

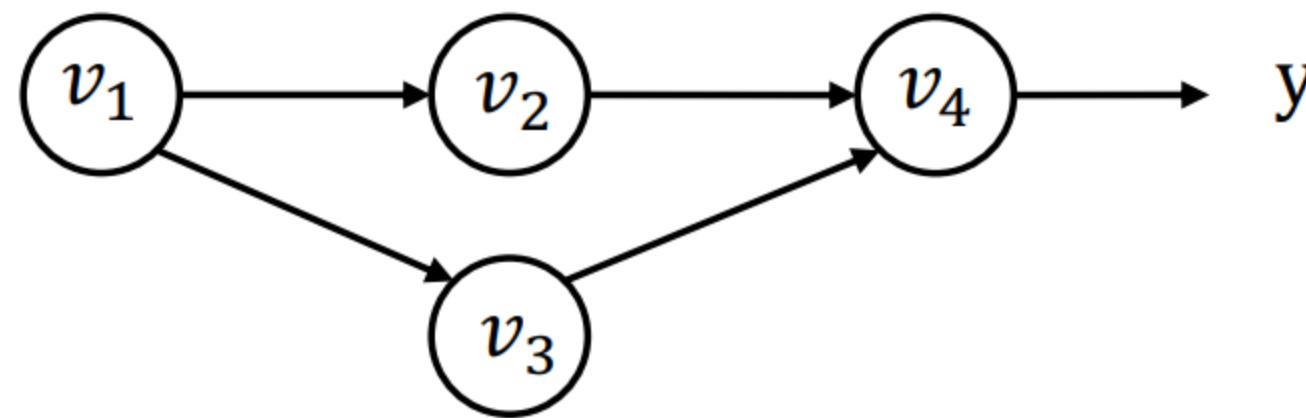
$$\bar{v}_3 = \bar{v}_6 \frac{\partial v_6}{\partial v_3} = \bar{v}_6 \times 1 = 1$$

$$\bar{v}_2 = \bar{v}_5 \frac{\partial v_5}{\partial v_2} + \bar{v}_4 \frac{\partial v_4}{\partial v_2} = \bar{v}_5 \times \cos v_2 + \bar{v}_4 \times v_1 = -0.284 + 2 = 1.716$$

$$\bar{v}_1 = \bar{v}_4 \frac{\partial v_4}{\partial v_1} + \bar{v}_3 \frac{\partial v_3}{\partial v_1} = \bar{v}_4 \times v_2 + \bar{v}_3 \frac{1}{v_1} = 5 + \frac{1}{2} = 5.5$$

- Finally:  $\frac{\partial y}{\partial x_1} = \bar{v}_1 = 5.5$

# Case Study



How to derive the gradient of  $v_1$

$$\bar{v}_1 = \frac{\partial y}{\partial v_1} = \frac{\partial f(v_2, v_3)}{\partial v_2} \frac{\partial v_2}{\partial v_1} + \frac{\partial f(v_2, v_3)}{\partial v_3} \frac{\partial v_3}{\partial v_1} \quad \bar{v}_2 = \bar{v}_2 \frac{\partial v_2}{\partial v_1} + \bar{v}_3 \frac{\partial v_3}{\partial v_1}$$

For a  $v_i$  used by multiple consumers:

$$\bar{v}_i = \sum_{j \in next(i)} \bar{v}_{i \rightarrow j} \quad , \text{ where } \bar{v}_{i \rightarrow j} = \bar{v}_j \frac{\partial v_j}{\partial v_i}$$

# Summary: Backward Mode Autodiff

- Start from the output nodes
- Derive gradient all the way back to the input node

# Back to Our Question

A repr that expresses the computation using primitives

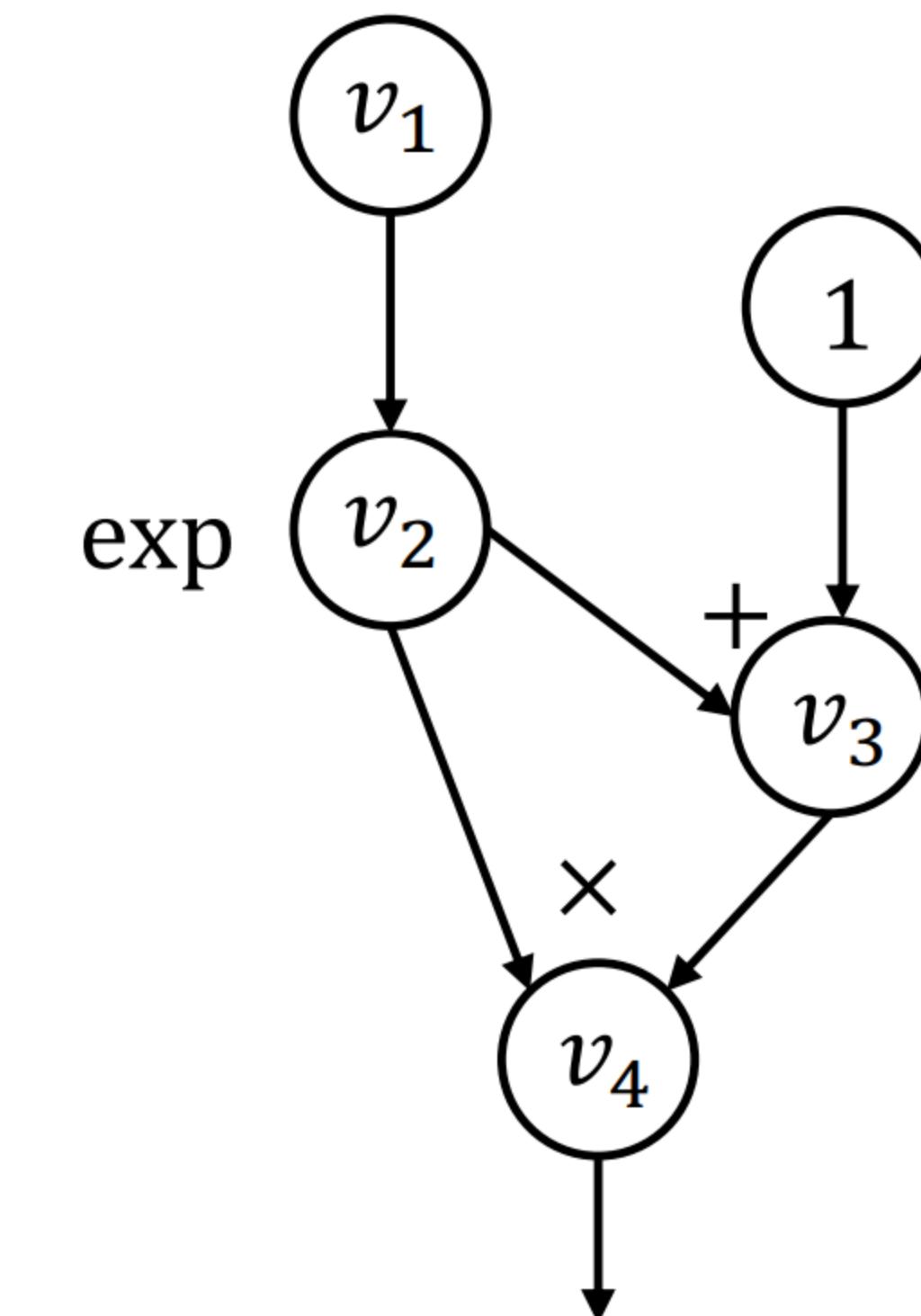
 A repr that expresses the **forward** computation using primitives

 A repr that expresses the **backward** computation using primitives

# Back to our question: Construct the Backward Graph

- How can we construct a computational graph that calculates the adjoint value?

```
def gradient(out):
    node_to_grad = {out: [1]}
    for i in reverse_topo_order(out):
         $\bar{v}_i = \sum_j \bar{v}_{i \rightarrow j} = \text{sum}(\text{node\_to\_grad}[i])$ 
        for k ∈ inputs(i):
            compute  $\bar{v}_{k \rightarrow i} = \bar{v}_i \frac{\partial v_i}{\partial v_k}$ 
            append  $\bar{v}_{k \rightarrow i}$  to node_to_grad[k]
    return adjoint of input  $\bar{v}_{input}$ 
```



$$f: (\exp(v_1) + 1)\exp(v_1)$$

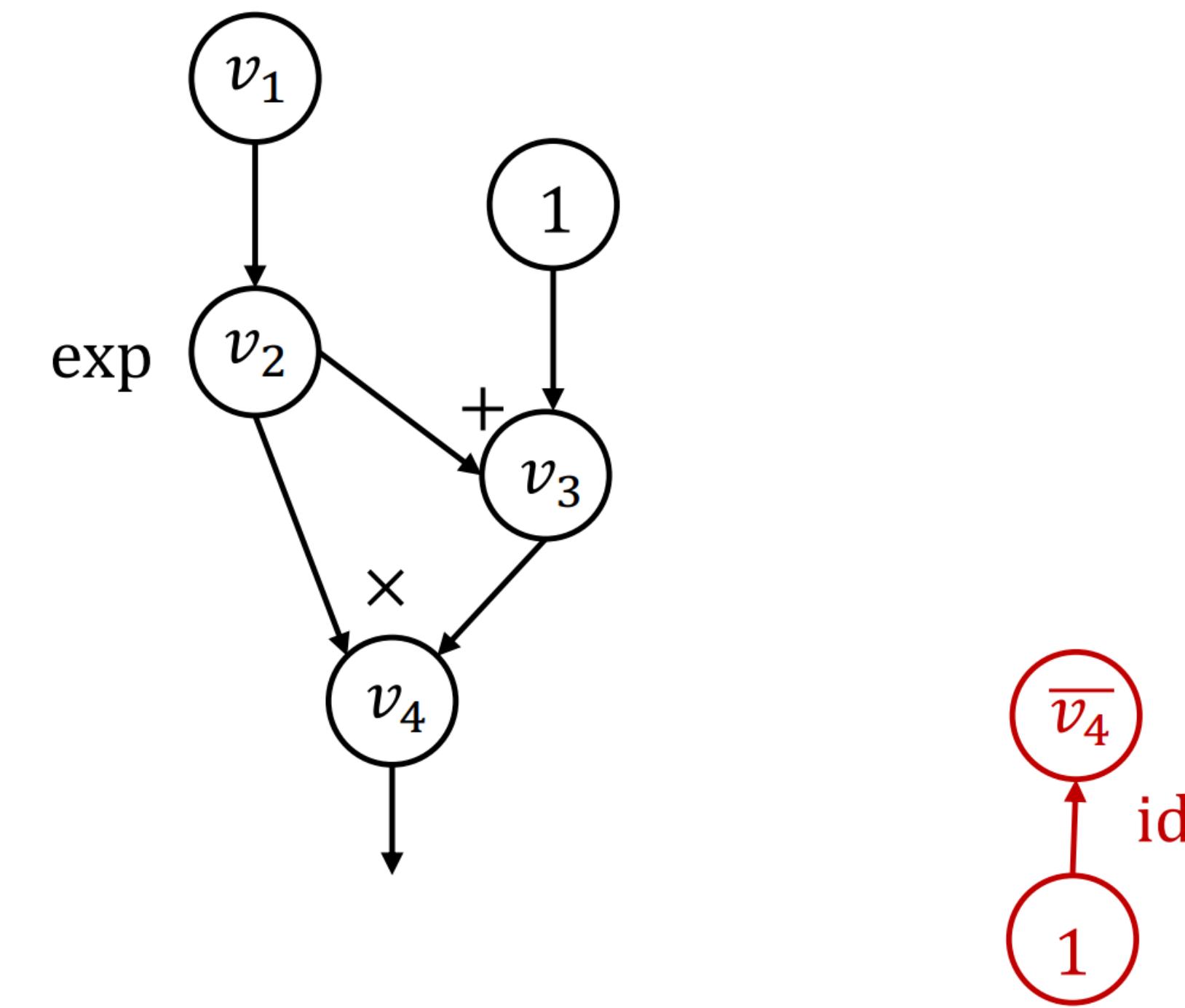
# How to implement reverse Autodiff (aka. BP)

Start from  $v_4$

$$i = 4: v_4 = \text{sum}([1]) = 1$$

```
def gradient(out):
    node_to_grad = {out: [1]}
    for i in reverse_topo_order(out):
         $\bar{v}_i = \sum_j \bar{v}_{i \rightarrow j} = \text{sum}(\text{node\_to\_grad}[i])$ 
        for k in inputs(i):
            compute  $\bar{v}_{k \rightarrow i} = \bar{v}_i \frac{\partial v_i}{\partial v_k}$ 
            append  $\bar{v}_{k \rightarrow i}$  to node_to_grad[k]
    return adjoint of input  $\bar{v}_{\text{input}}$ 
```

$i = 4$   
node\_to\_grad: {  
  4: [ $\bar{v}_4$ ]  
}

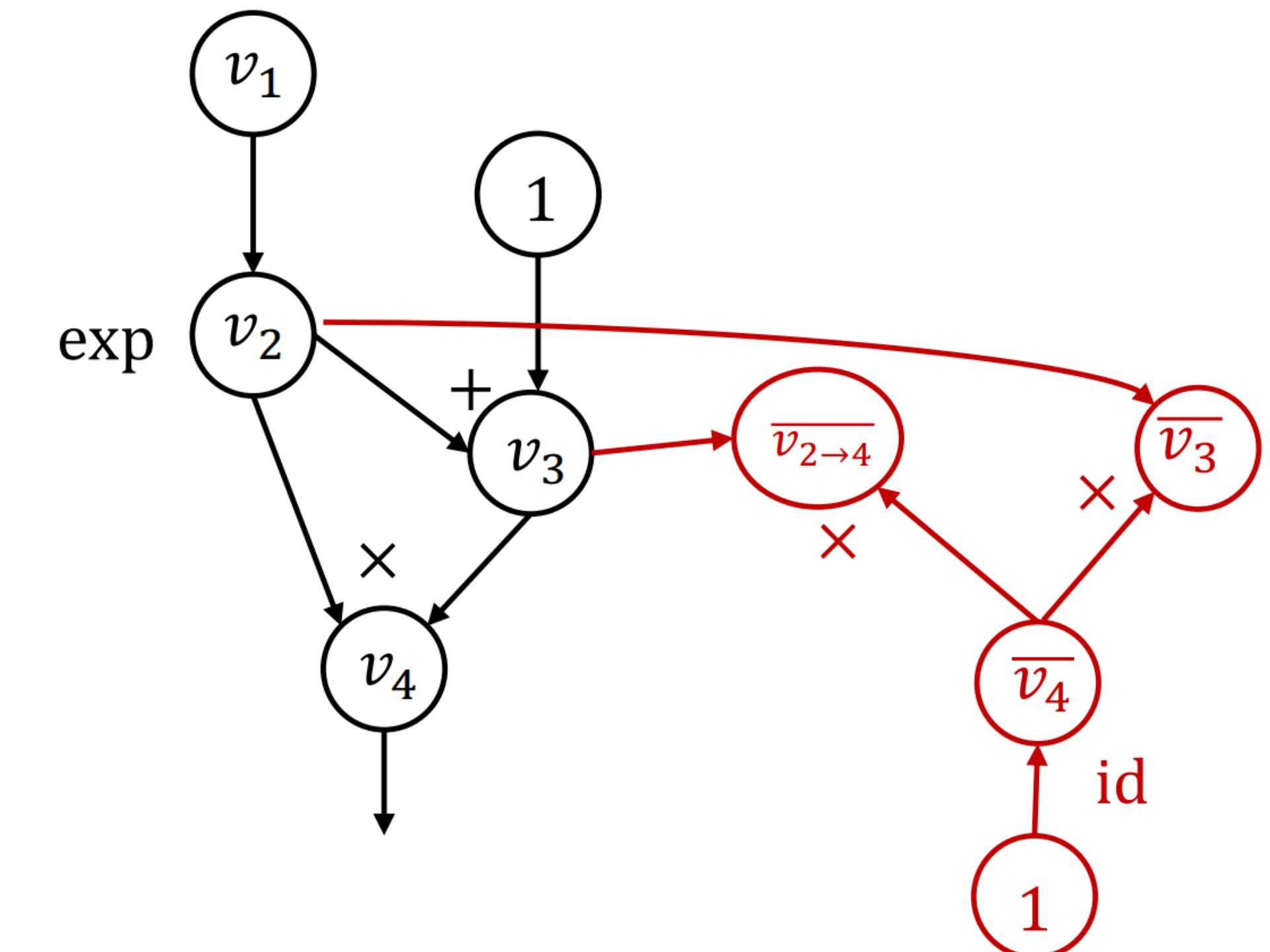


$v_4$ : Inspect  $(v_2, v_4)$  and  $(v_3, v_4)$

```
def gradient(out):
    node_to_grad = {out: [1]}
    for i in reverse_topo_order(out):
         $\bar{v}_i = \sum_j \bar{v}_{i \rightarrow j} = \text{sum}(\text{node\_to\_grad}[i])$ 
        for k in inputs(i):
            compute  $\bar{v}_{k \rightarrow i} = \bar{v}_i \frac{\partial v_i}{\partial v_k}$ 
            append  $\bar{v}_{k \rightarrow i}$  to node_to_grad[k]
    return adjoint of input  $\bar{v}_{\text{input}}$ 
```

$i = 4$   
 $\text{node\_to\_grad: } \{$   
 $2: [\bar{v}_{2 \rightarrow 4}]$   
 $3: [\bar{v}_3]$   
 $4: [\bar{v}_4]$   
 $\}$

$$\begin{aligned} i=4: \bar{v}_4 &= \text{sum}([1]) = 1 \\ k=2: \bar{v}_{2 \rightarrow 4} &= \bar{v}_4 \frac{\partial v_4}{\partial v_2} = \bar{v}_4 v_3 \\ k=3: \bar{v}_{3 \rightarrow 4} &= \bar{v}_4 \frac{\partial v_4}{\partial v_3} = \bar{v}_4 v_2, \bar{v}_{3 \rightarrow 4} = \bar{v}_3 \end{aligned}$$



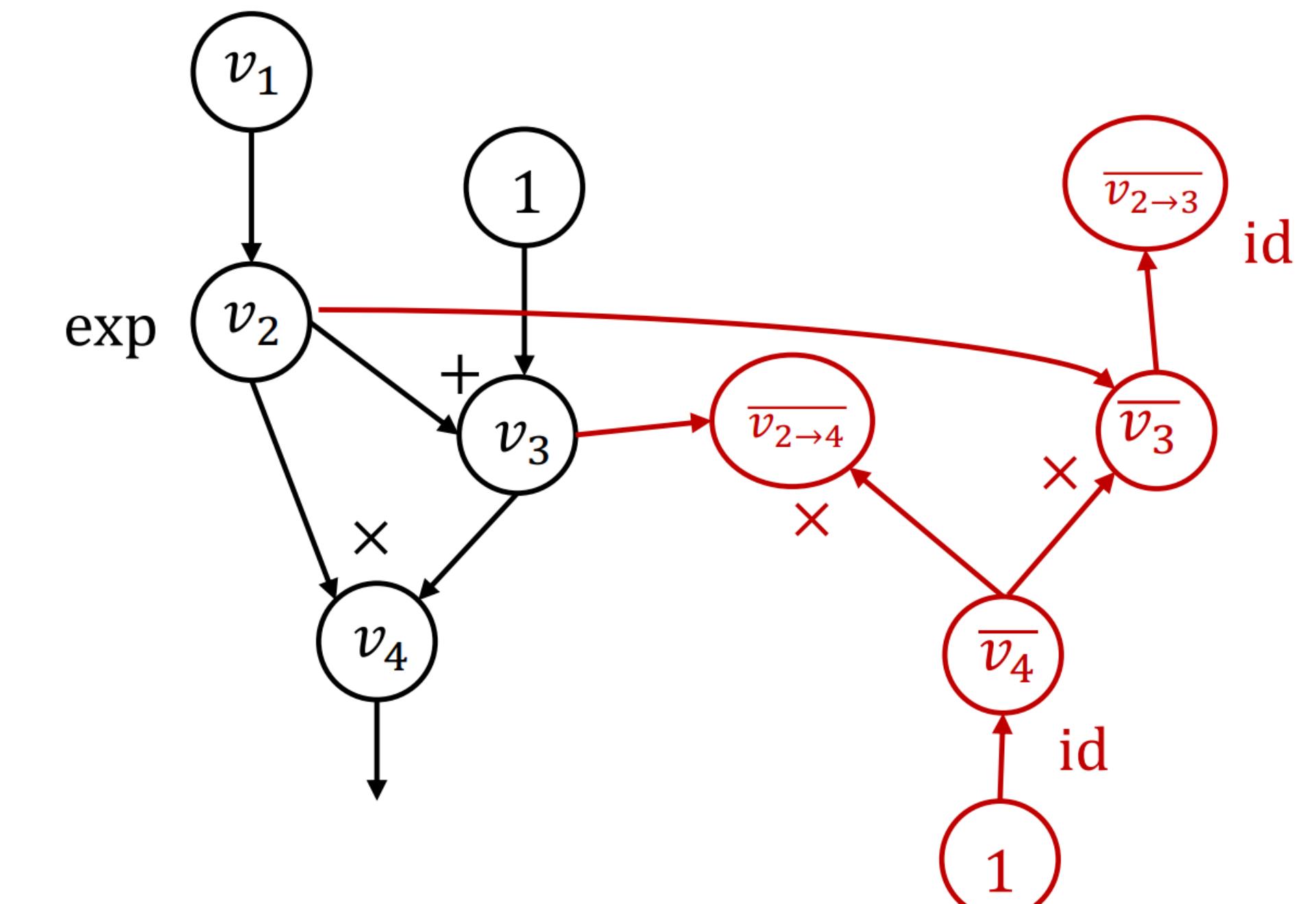
Inspect  $v_3$

```
def gradient(out):
    node_to_grad = {out: [1]}
    for i in reverse_topo_order(out):
         $\bar{v}_i = \sum_j \bar{v}_{i \rightarrow j} = \text{sum}(\text{node\_to\_grad}[i])$ 
        for k ∈ inputs(i):
            compute  $\bar{v}_{k \rightarrow i} = \bar{v}_i \frac{\partial v_i}{\partial v_k}$ 
            append  $\bar{v}_{k \rightarrow i}$  to node_to_grad[k]
    return adjoint of input  $\bar{v}_{input}$ 
```

$i = 3$   
 $\text{node\_to\_grad: } \{$   
 $2: [\bar{v}_{2 \rightarrow 4}, \bar{v}_{2 \rightarrow 3}]$   
 $3: [\bar{v}_3]$   
 $4: [\bar{v}_4]$   
 $\}$

$i=3: \bar{v}_3$  done!

$$k=2: \bar{v}_{2 \rightarrow 3} = \bar{v}_3 \frac{\partial v_3}{\partial v_2} = \bar{v}_3$$

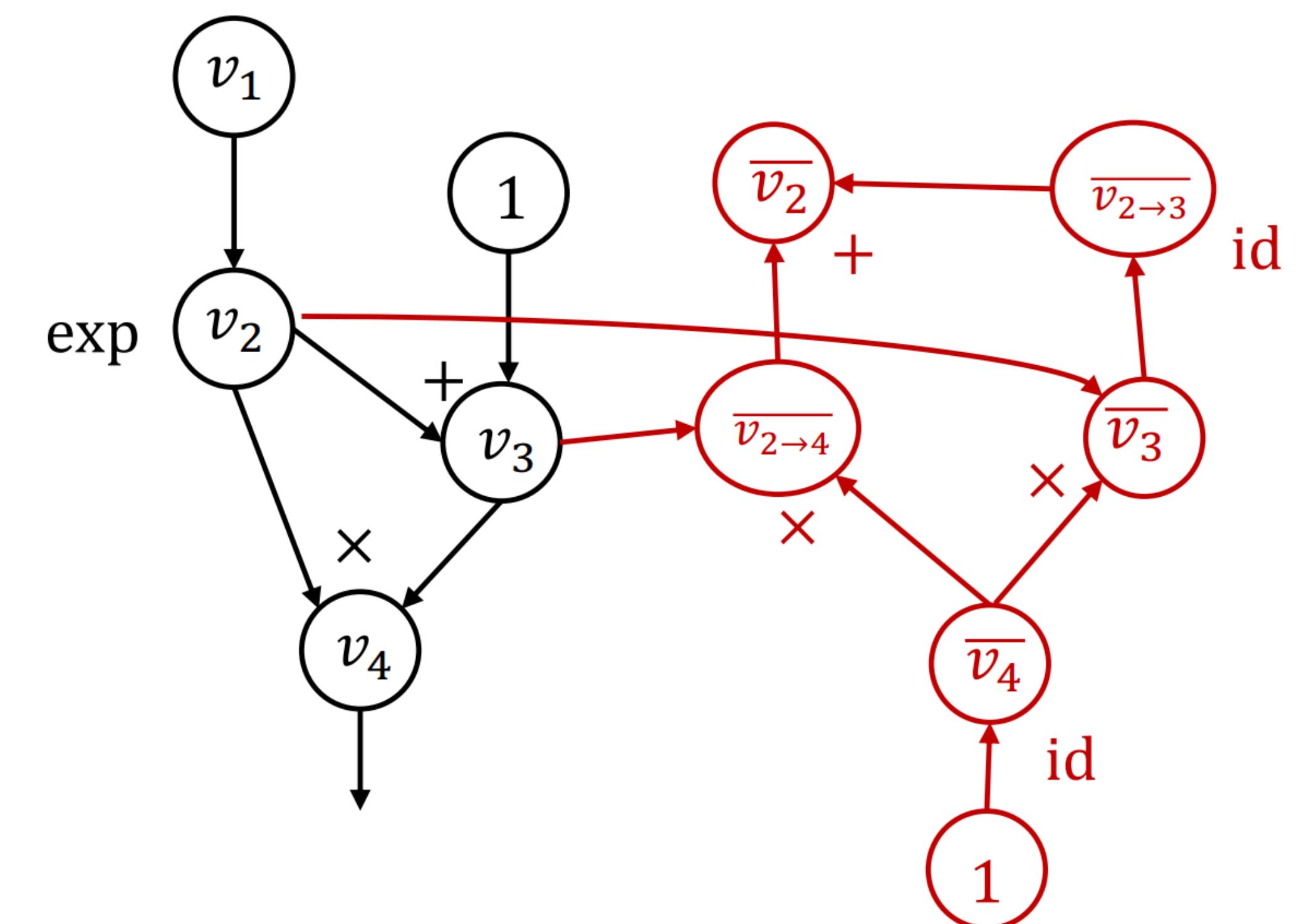


$$i=2: \bar{v}_2 = \overline{v_{2 \rightarrow 3}} + \overline{v_{2 \rightarrow 4}}$$

Inspect  $v_2$

```
def gradient(out):
    node_to_grad = {out: [1]}
    for i in reverse_topo_order(out):
         $\bar{v}_i = \sum_j \bar{v}_{i \rightarrow j} = \text{sum}(\text{node\_to\_grad}[i])$ 
        for k in inputs(i):
            compute  $\bar{v}_{k \rightarrow i} = \bar{v}_i \frac{\partial v_i}{\partial v_k}$ 
            append  $\bar{v}_{k \rightarrow i}$  to node_to_grad[k]
    return adjoint of input  $\bar{v}_{input}$ 
```

$i = 2$   
 $\text{node\_to\_grad: } \{$   
 $2: [\bar{v}_{2 \rightarrow 4}, \bar{v}_{2 \rightarrow 3}]$   
 $3: [\bar{v}_3]$   
 $4: [\bar{v}_4]$   
 $\}$



# Inspect $(v_1, v_2)$

```

def gradient(out):
    node_to_grad = {out: [1]}
    for i in reverse_topo_order(out):
         $\bar{v}_i = \sum_j \bar{v}_{i \rightarrow j} = \text{sum}(\text{node\_to\_grad}[i])$ 
        for k in inputs(i):
            compute  $\bar{v}_{k \rightarrow i} = \bar{v}_i \frac{\partial v_i}{\partial v_k}$ 
            append  $\bar{v}_{k \rightarrow i}$  to node_to_grad[k]
    return adjoint of input  $\bar{v}_{input}$ 

```



$i = 2$

```

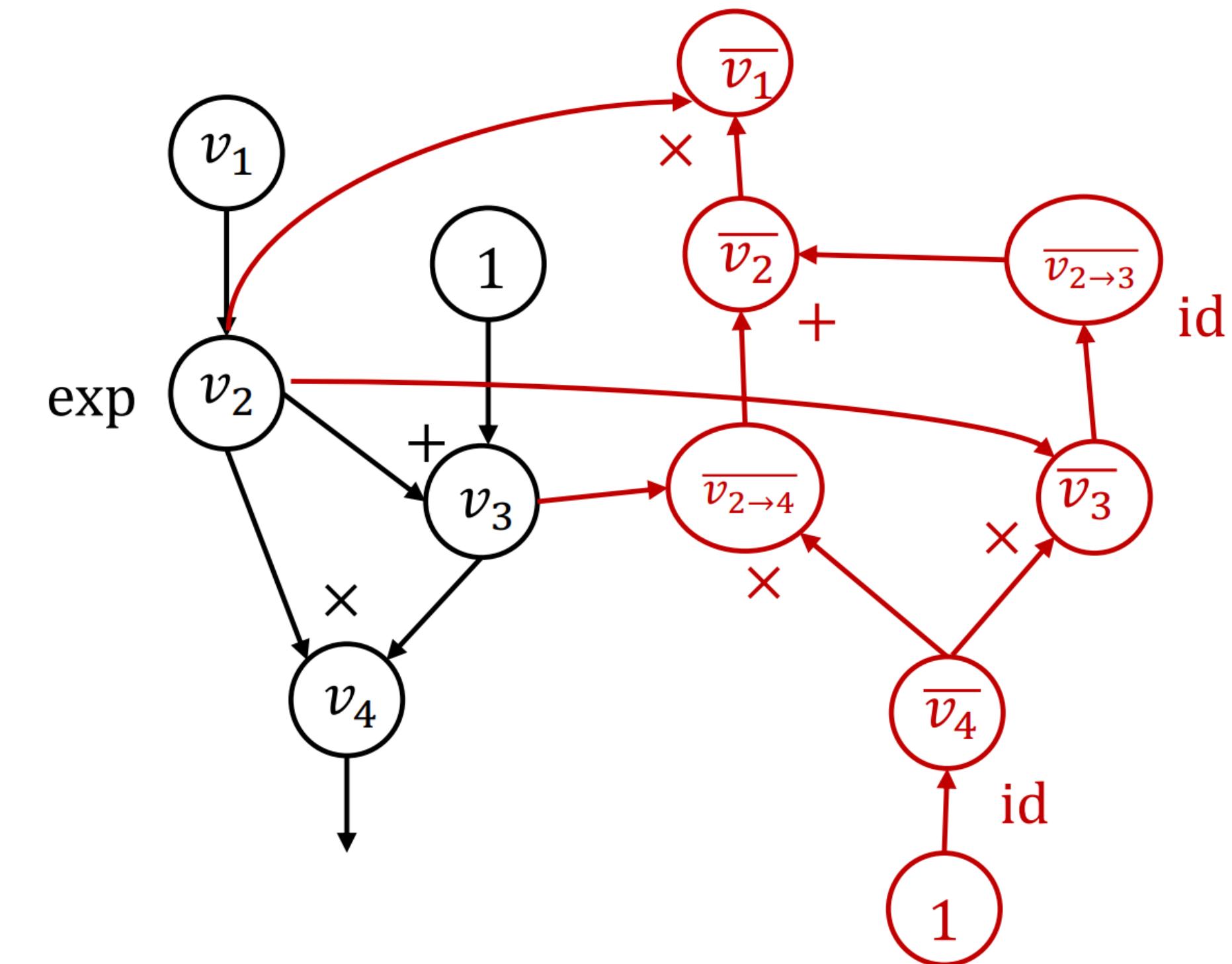
node_to_grad: {
    1: [ $\bar{v}_1$ ]
    2: [ $\bar{v}_{2 \rightarrow 4}$ ,  $\bar{v}_{2 \rightarrow 3}$ ]
    3: [ $\bar{v}_3$ ]
    4: [ $\bar{v}_4$ ]
}

```

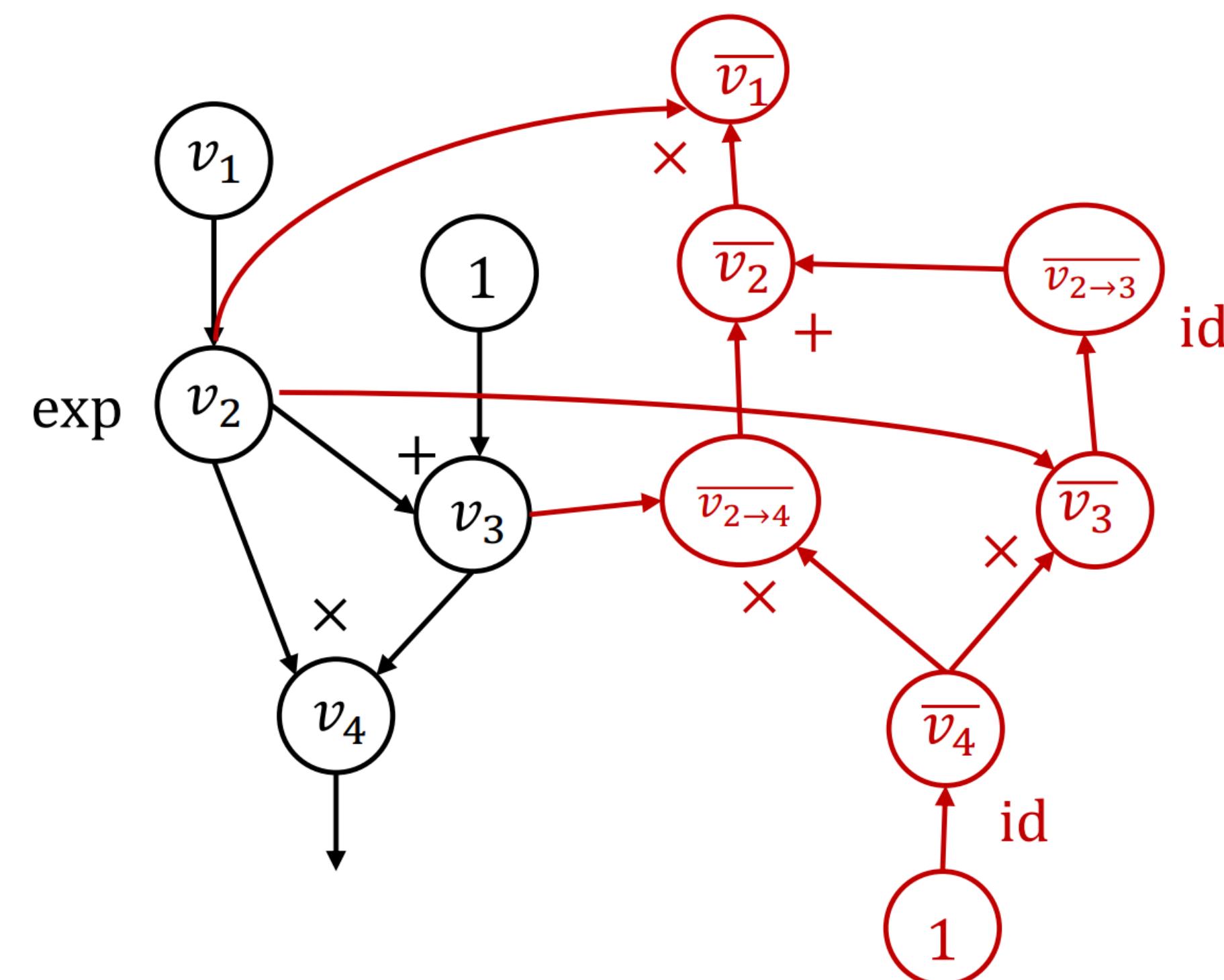
$$i=2: \bar{v}_2 = \bar{v}_{2 \rightarrow 3} + \bar{v}_{2 \rightarrow 4}$$

$$k=1: \bar{v}_{1 \rightarrow 2} = \bar{v}_2 \frac{\partial v_2}{\partial v_1} = \bar{v}_2 \exp(v_1),$$

$$\bar{v}_1 = \bar{v}_{1 \rightarrow 2}$$

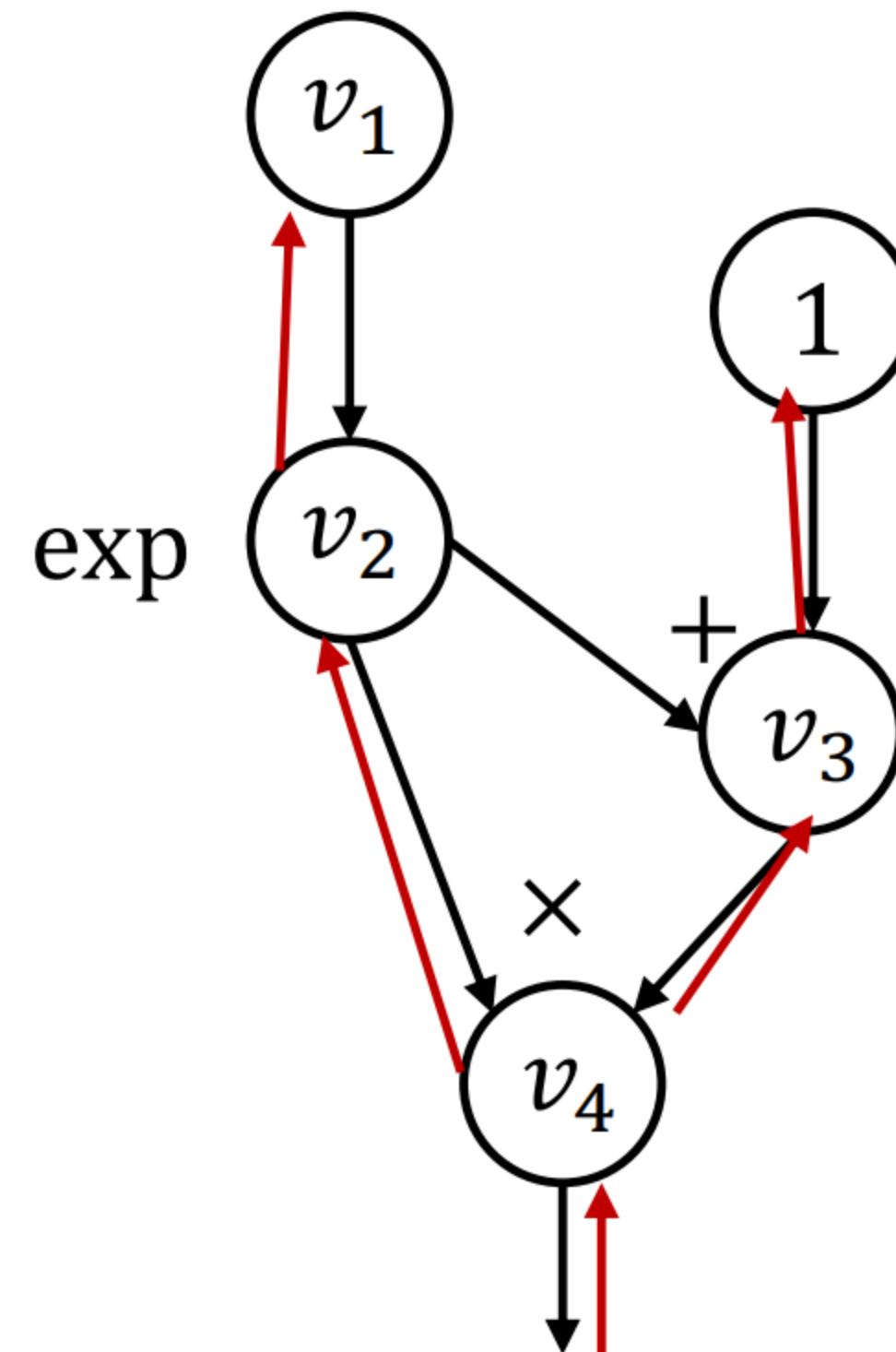


# Summary: Backward AD

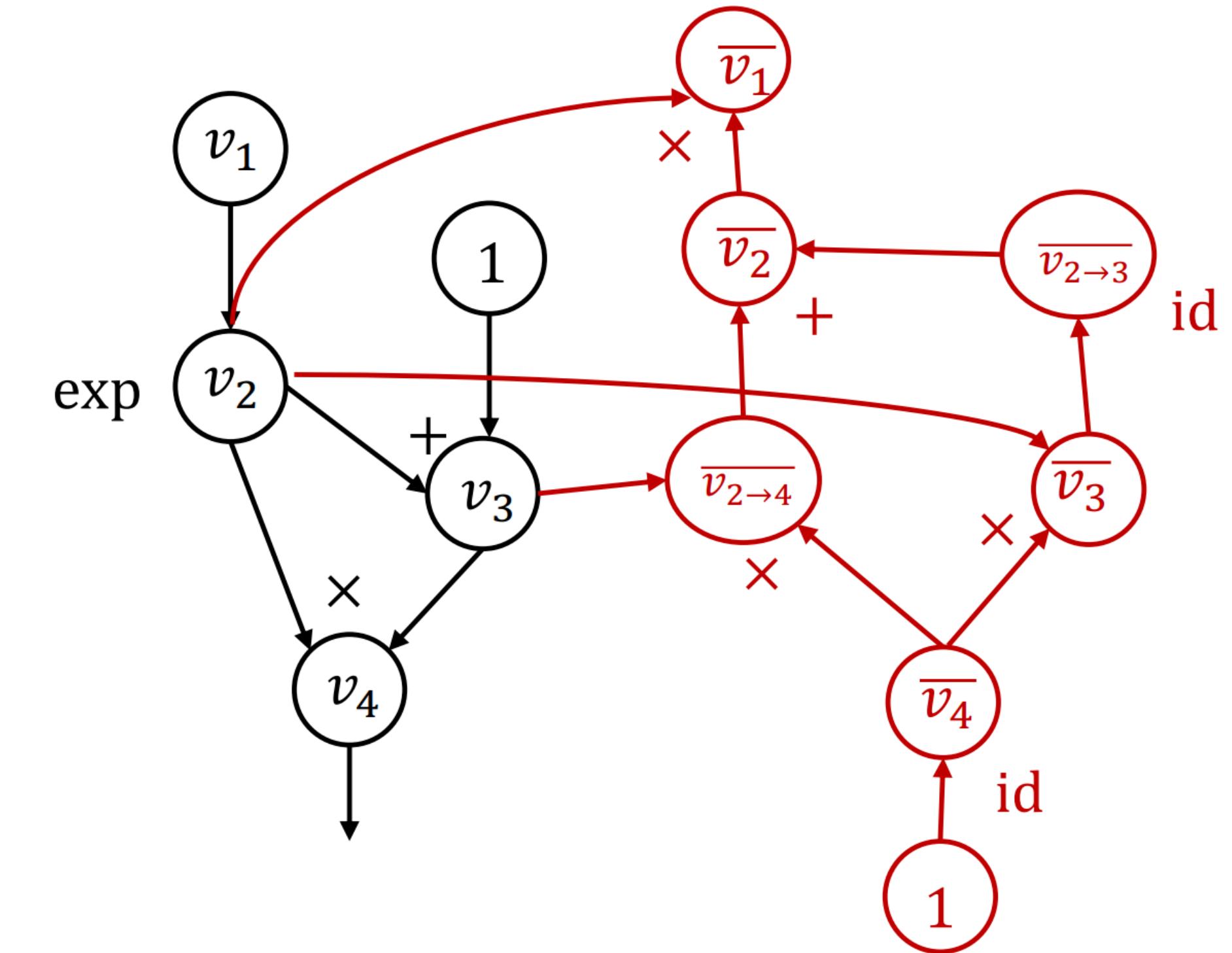


- Construct backward graph in a symbolic way (instead of concrete values)
- This graph can be reused by different input values

# Backpropagation vs. Reverse-mode AD



vs.

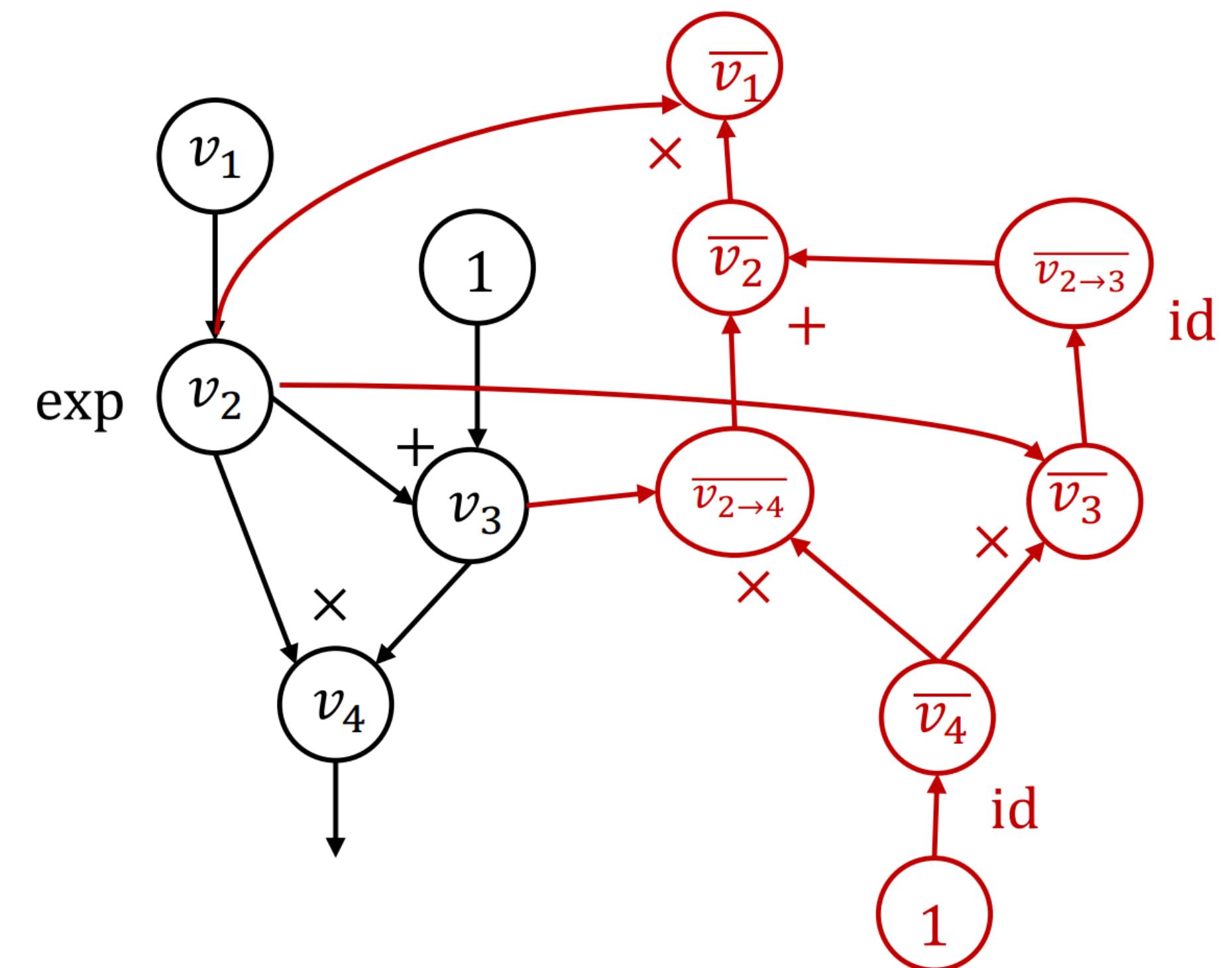


- Run backward through the forward graph
- Caffe/cuda-convnet

- Construct backward graph
- Used by TensorFlow, PyTorch

# Incomplete yet?

- What is the missing from the following graph for ML training?



# Recall Our Master Equation

$$\theta^{(t+1)} = f(\theta^{(t)}, \nabla_L(\theta^{(t)}, D^{(t)}))$$

$$L = \text{MSE}(w_2 \cdot \text{ReLU}(w_1 x), y) \quad \theta = \{w_1, w_2\}, D = \{(x, y)\}$$

$$f(\theta, \nabla_L) = \theta - \nabla_L$$

Forward

$$L(\cdot)$$

Backward

$$\nabla_L(\cdot)$$

Weight update

$$f(\cdot)$$

# Put in Practice

$$\theta^{(t+1)} = f(\theta^{(t)}, \nabla_L(\theta^{(t)}, D^{(t)}))$$

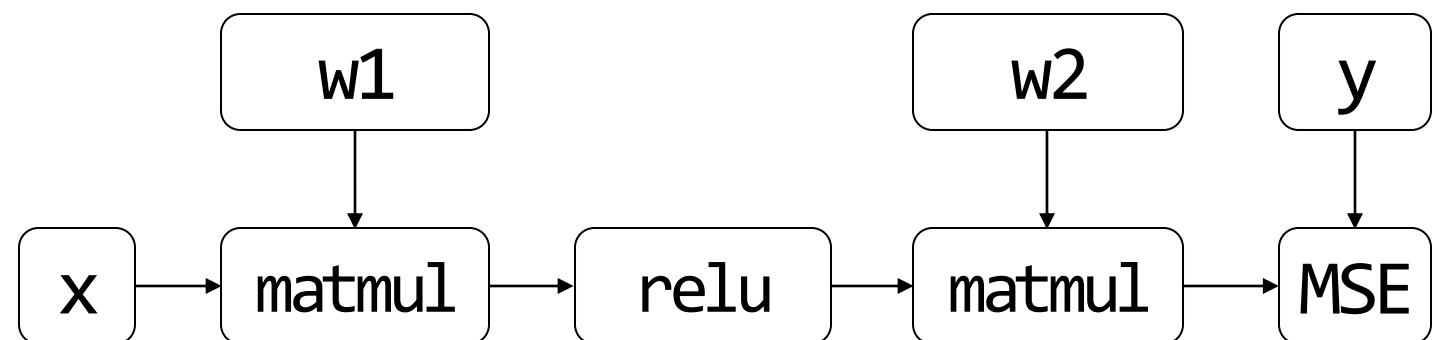
$$L = \text{MSE}(w_2 \cdot \text{ReLU}(w_1 x), y) \quad \theta = \{w_1, w_2\}, D = \{(x, y)\}$$

$$f(\theta, \nabla_L) = \theta - \nabla_L$$

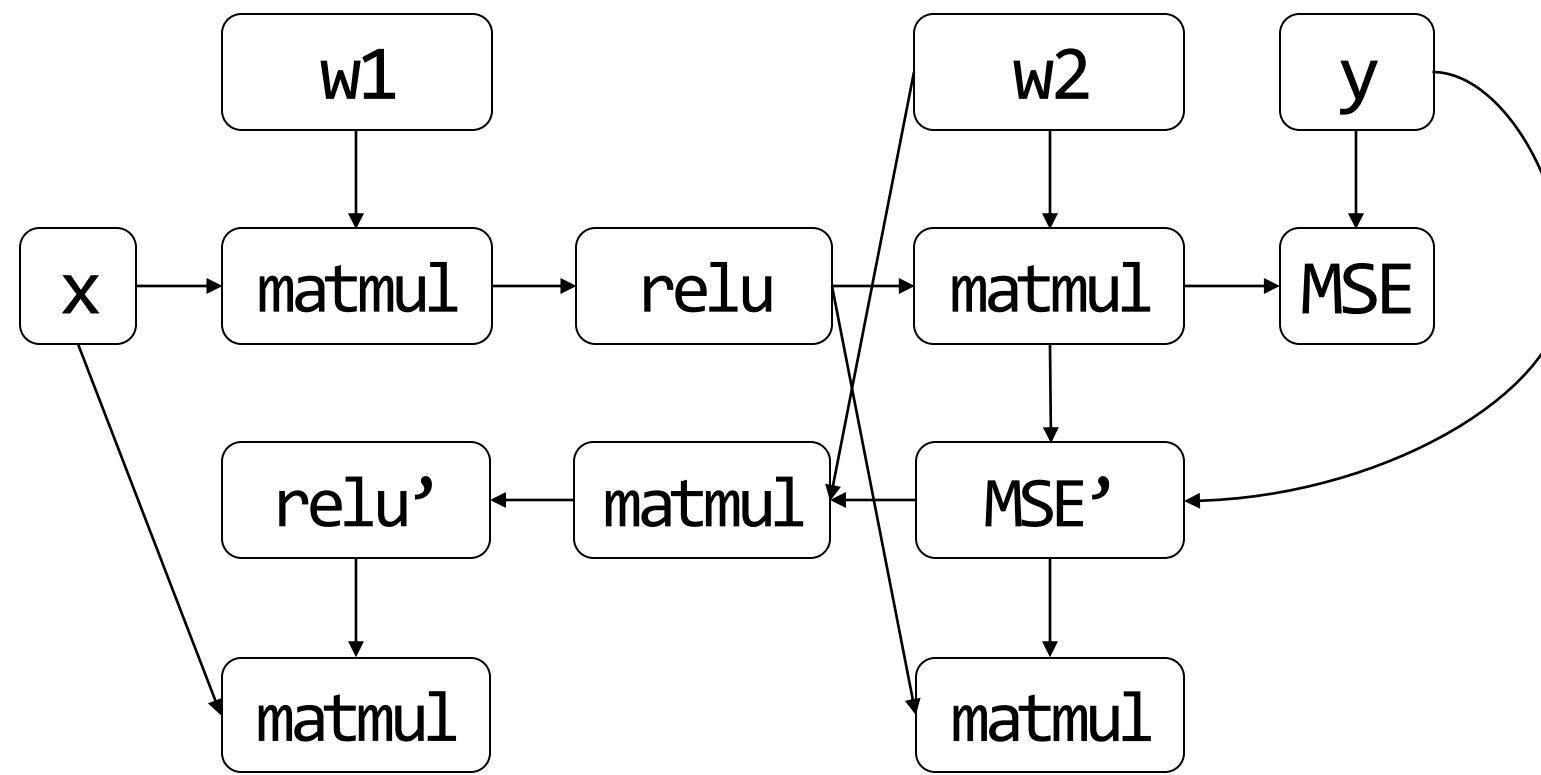
Operator / its output tensor

→ Data flowing direction

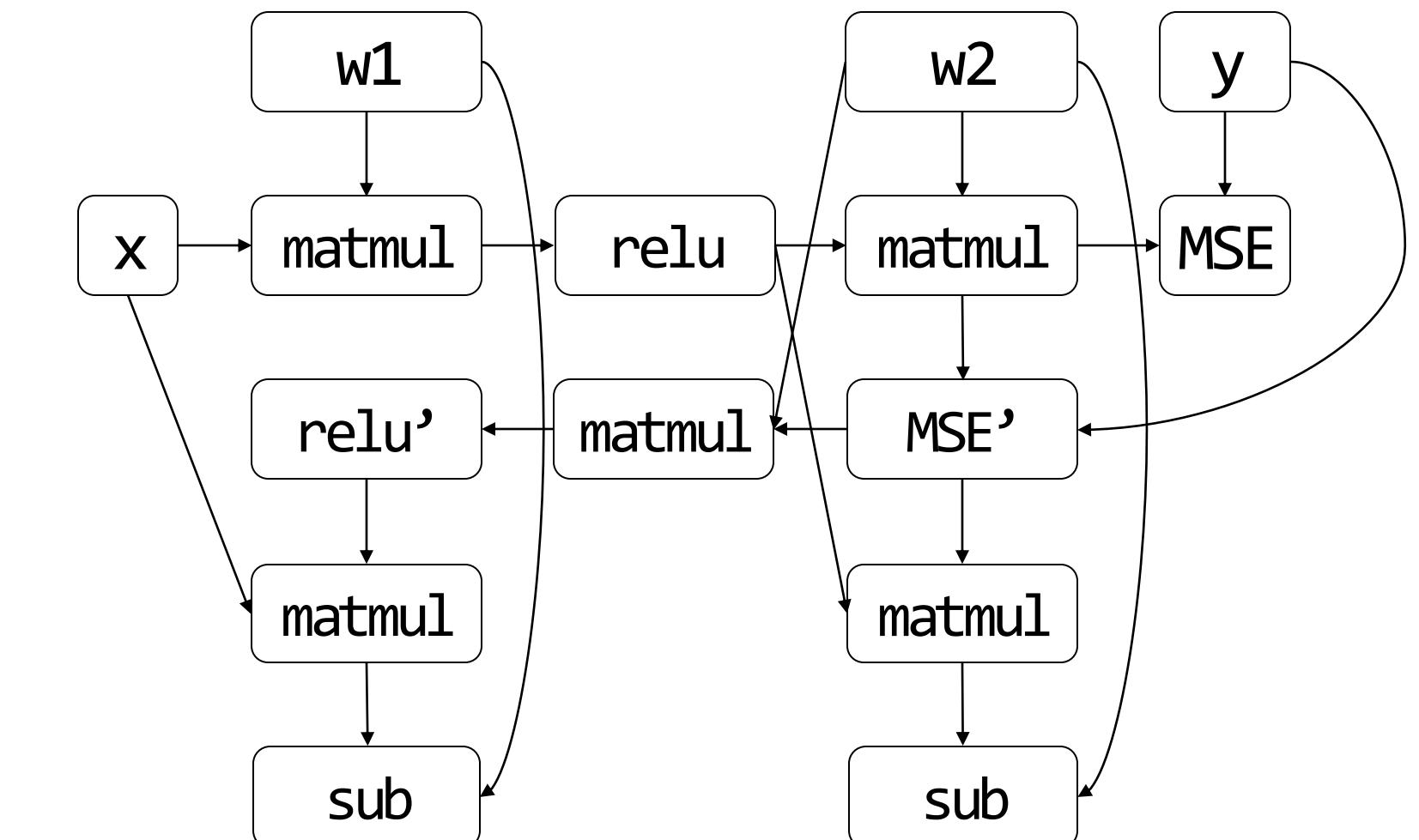
Forward



+Backward



+Weight update



# Homework: How to derive gradients for

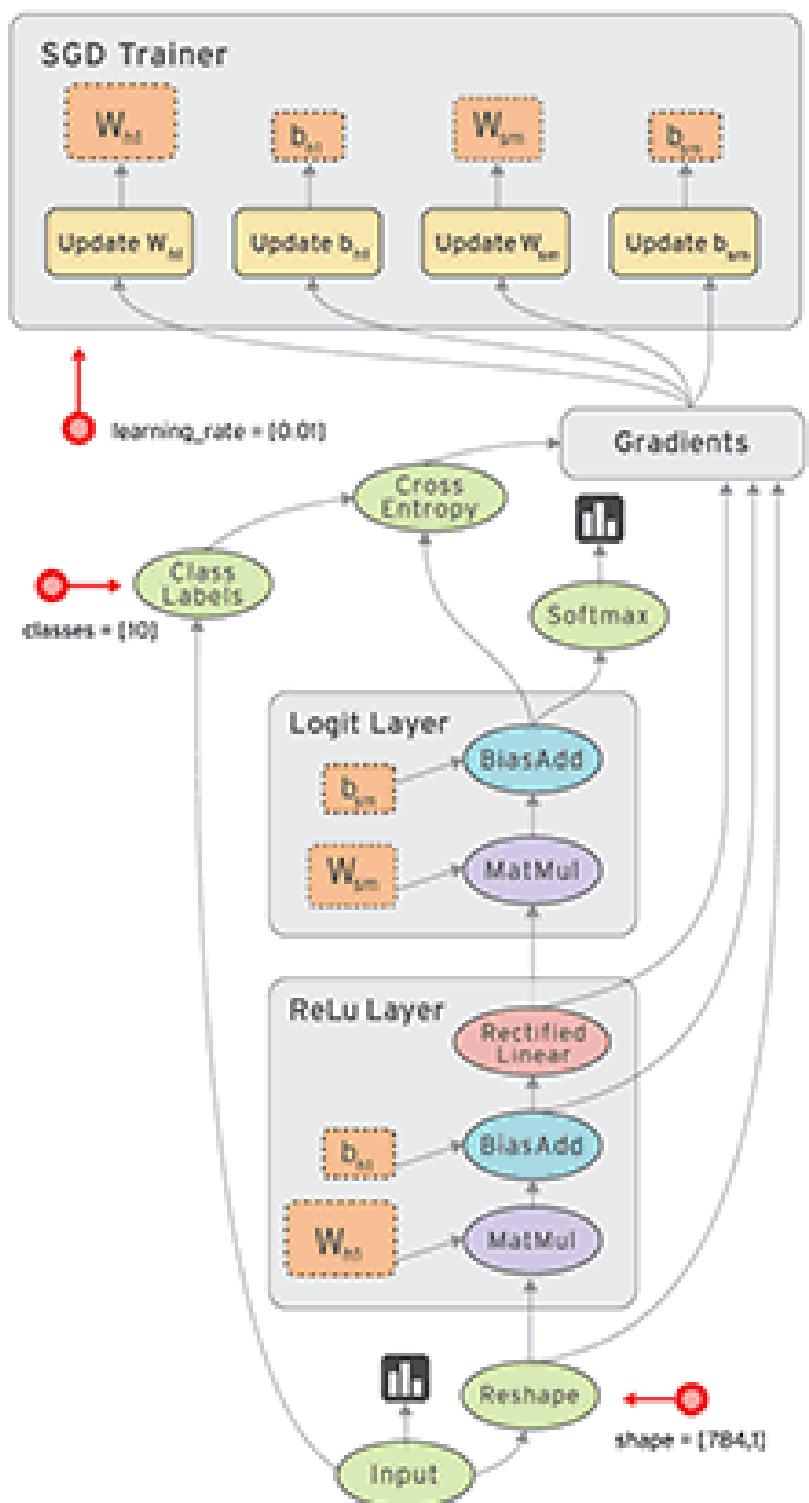
- Softmax cross entropy:

$$L = -\sum t_i \log(y_i), y_i = \text{softmax}(\mathbf{x})_i = \frac{e^{x_i}}{\sum e^{x_d}}$$

# Today

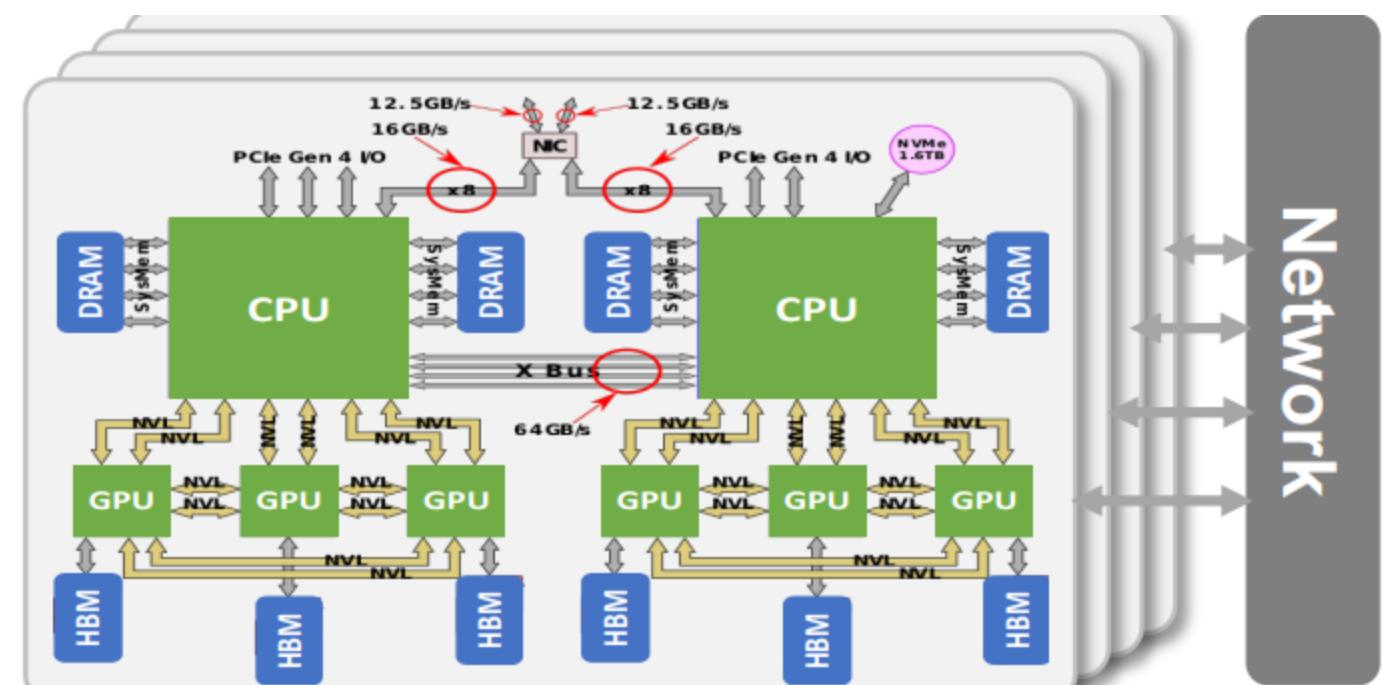
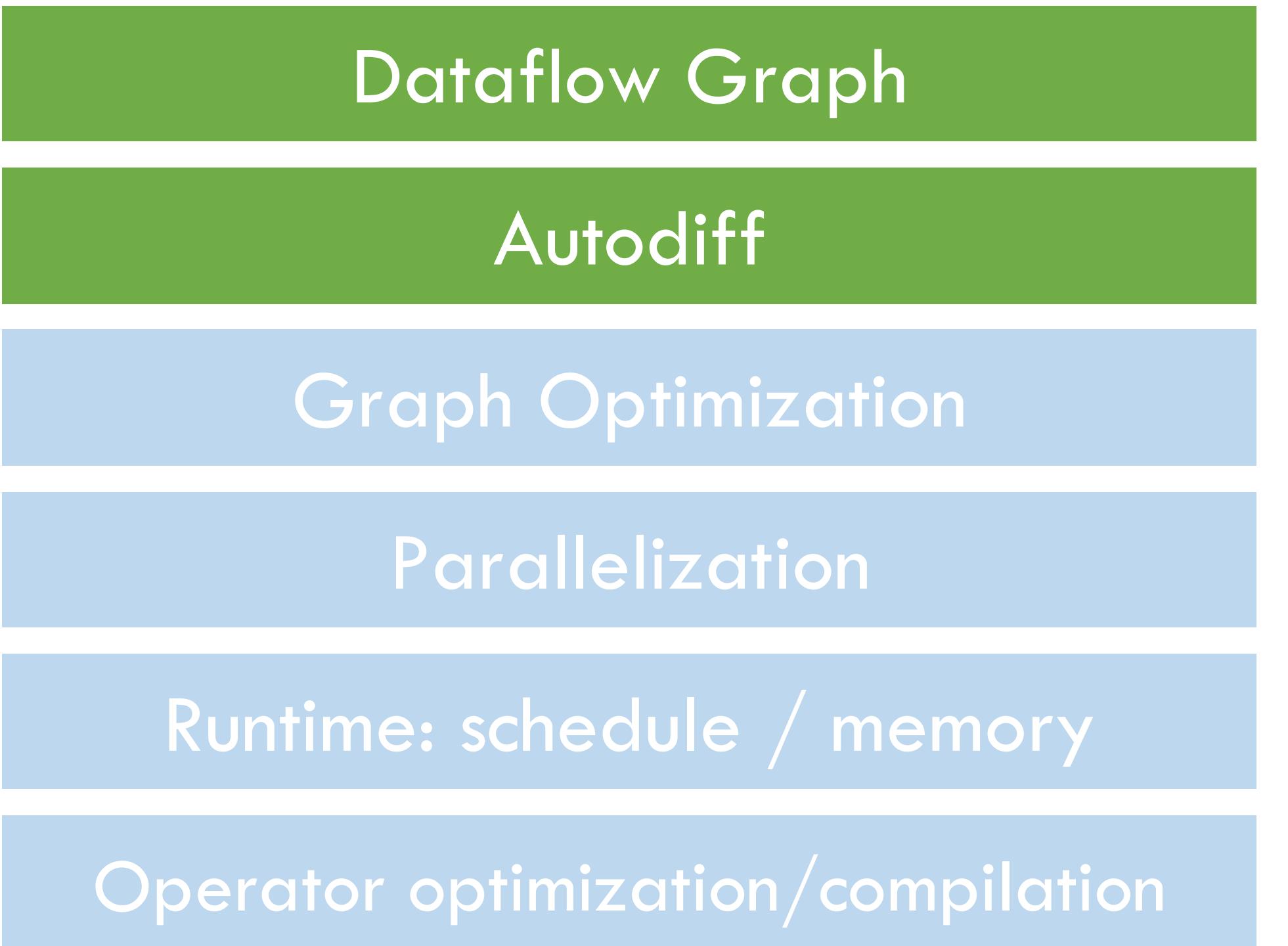
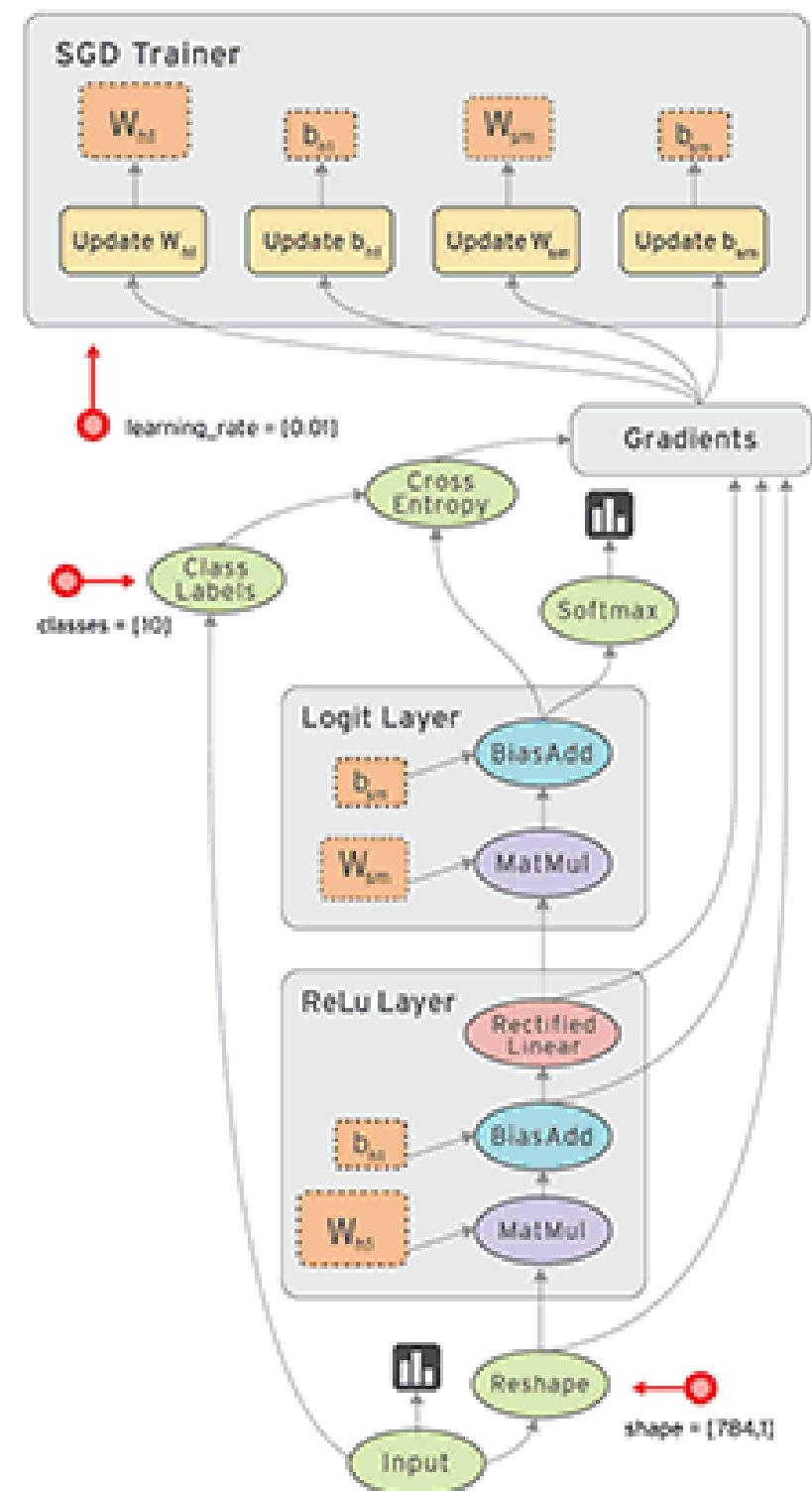
- Autodiff
- **Architecture Overview**

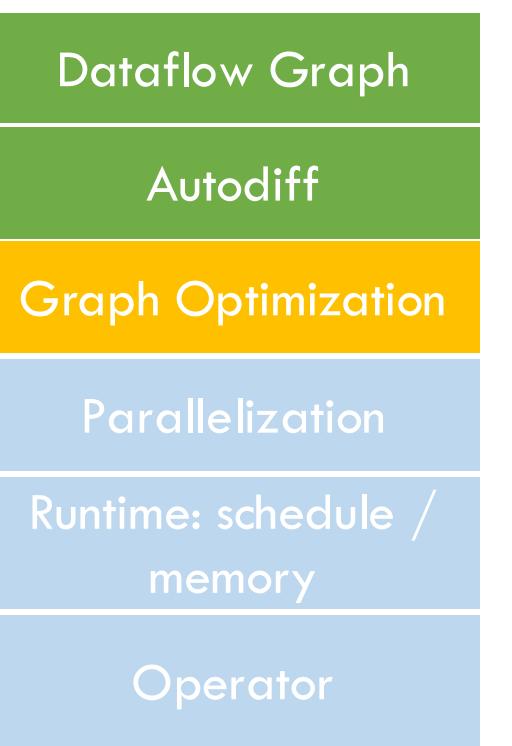
# MLSys' Grand problem



- Our system goals:
  - Fast
  - Scale
  - Memory-efficient
  - Run on diverse hardware
  - Energy-efficient
  - Easy to program/debug/deploy

# ML System Overview



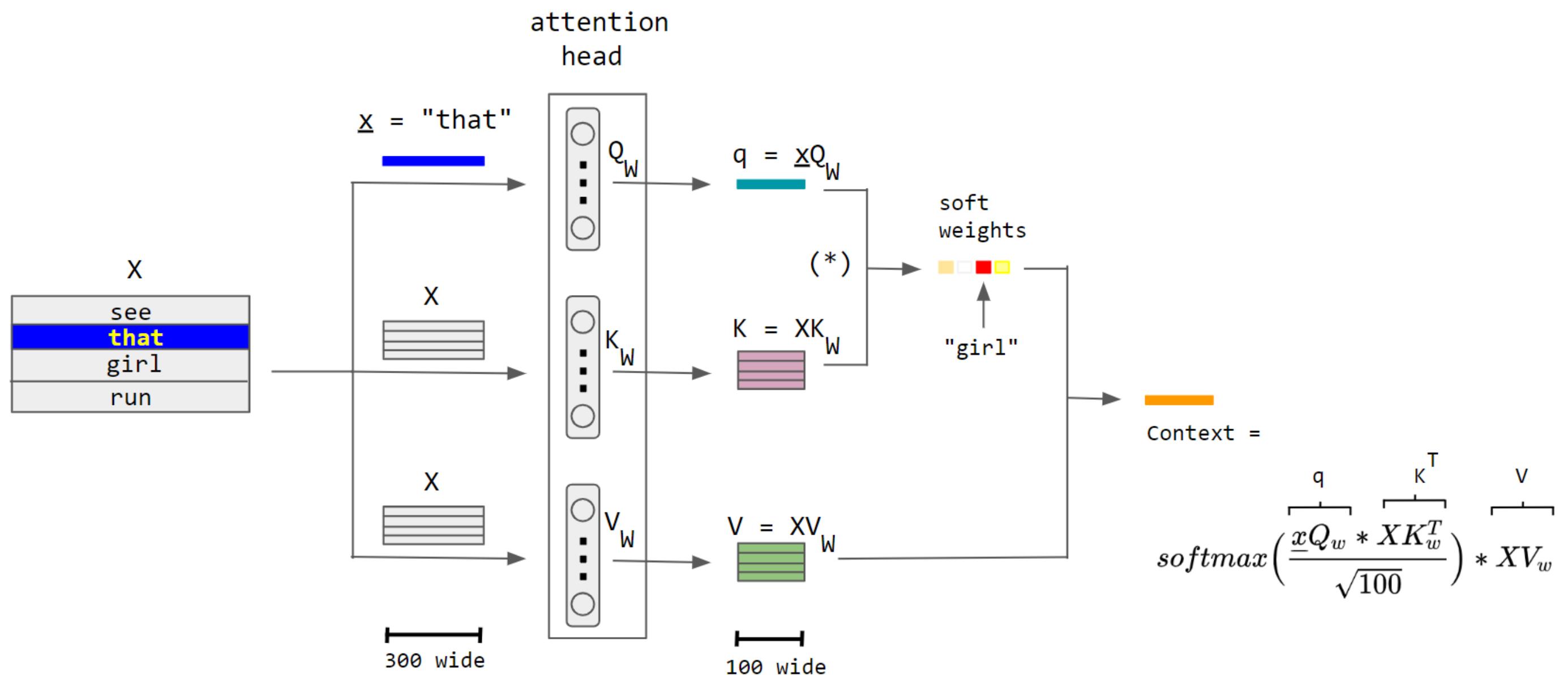


# Graph Optimization

- Goal:
  - Rewrite the original Graph  $G$  to  $G'$
  - $G'$  runs faster than  $G$

Dataflow Graph
Autodiff
Graph Optimization
Parallelization
Runtime: schedule / memory
Operator

# Motivating Example: Attention



# Original

```
Q = matmul(w_q, h)
K = matmul(w_k, h)
V = matmul(w_v, h)
```

# Merged QKV

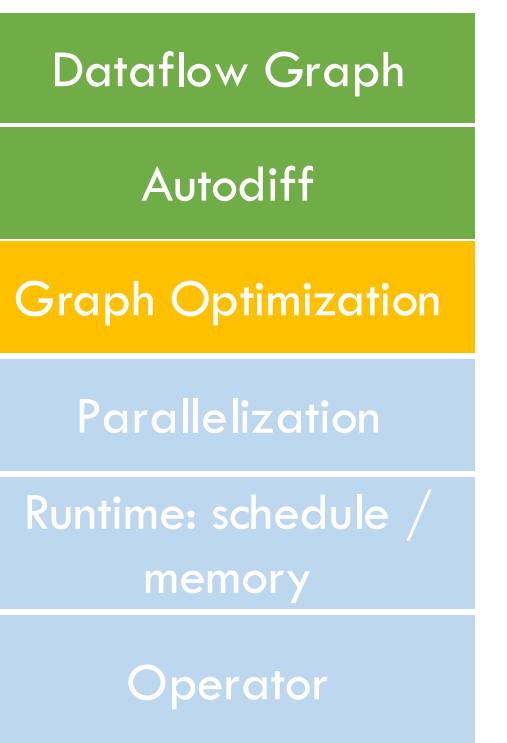
```
QKV = matmul(concat(w_q, w_k, w_v), h)
```

$$\text{softmax}\left(\frac{xQ_w * XK_w^T}{\sqrt{100}}\right) * XV_w$$

- Why merged QKV is faster?

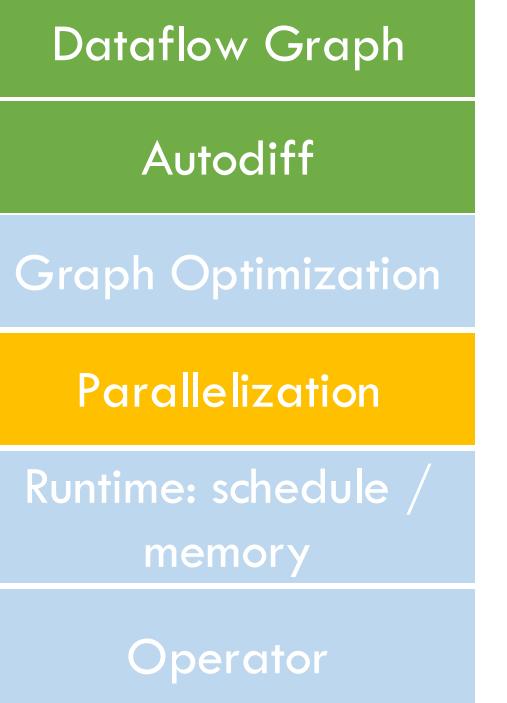
# Arithmetic Intensity

$$AI = \#ops / \#bytes$$



# How to perform graph optimization?

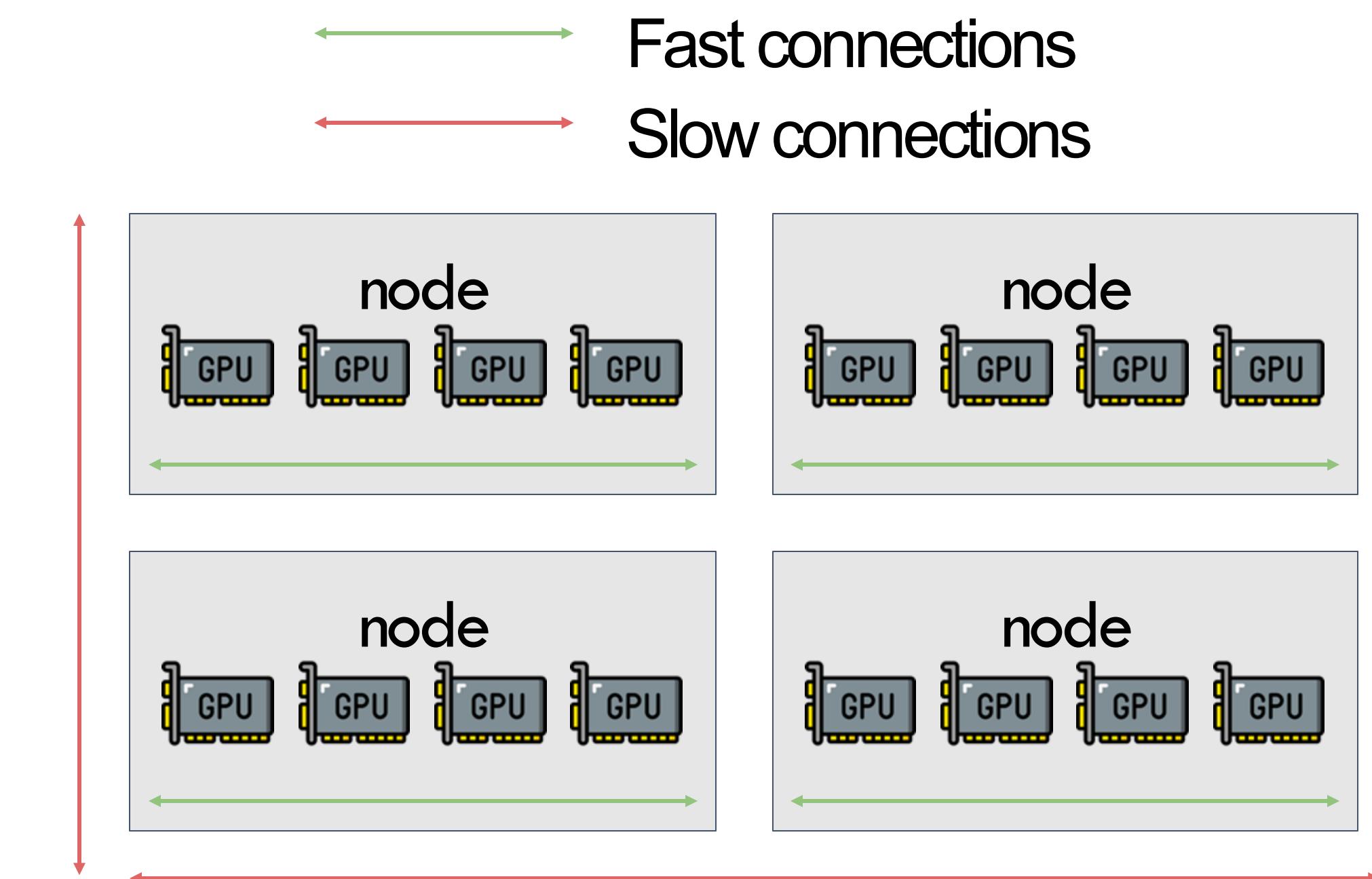
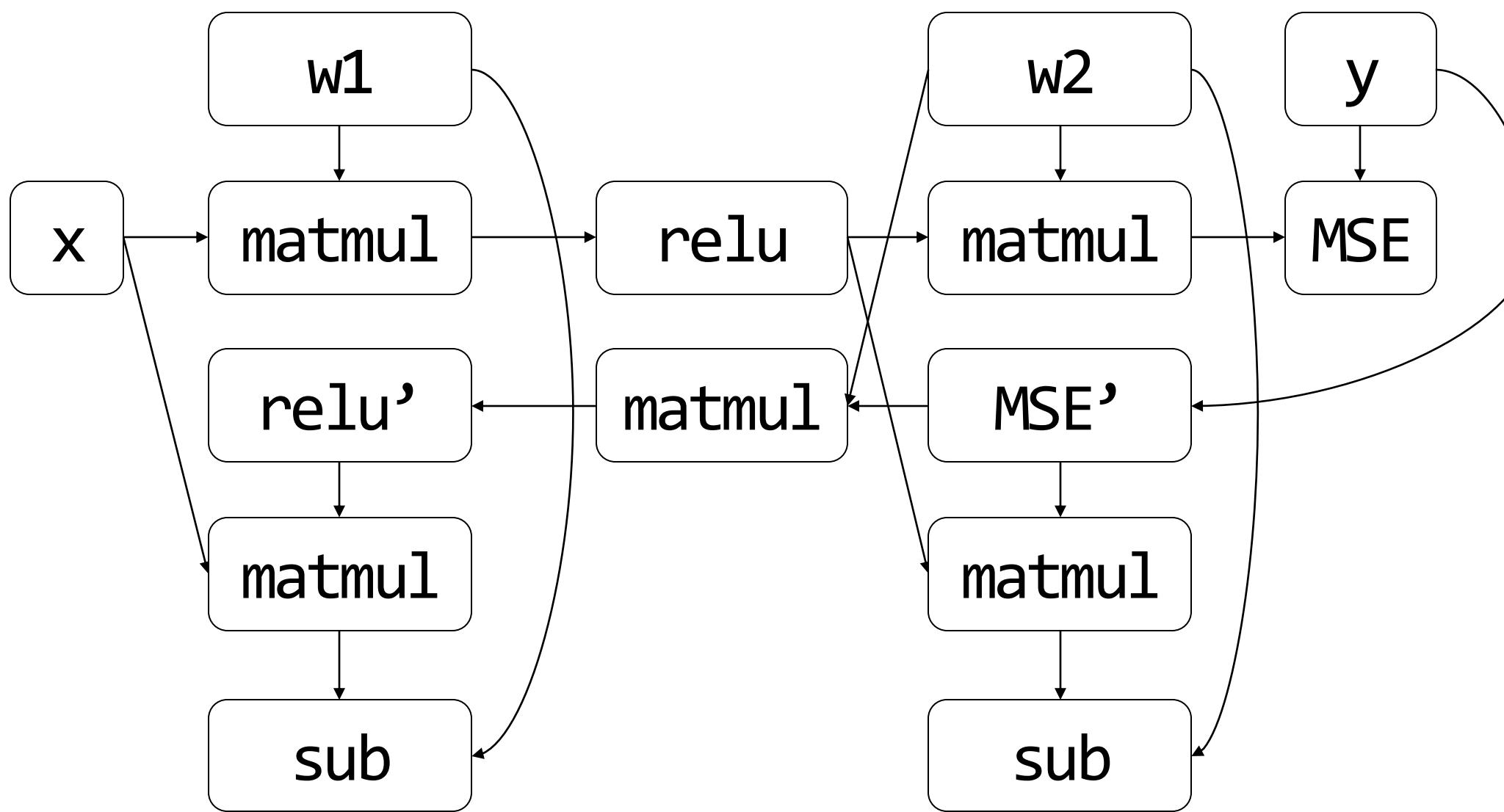
- Writing rules / template
- Auto discovery

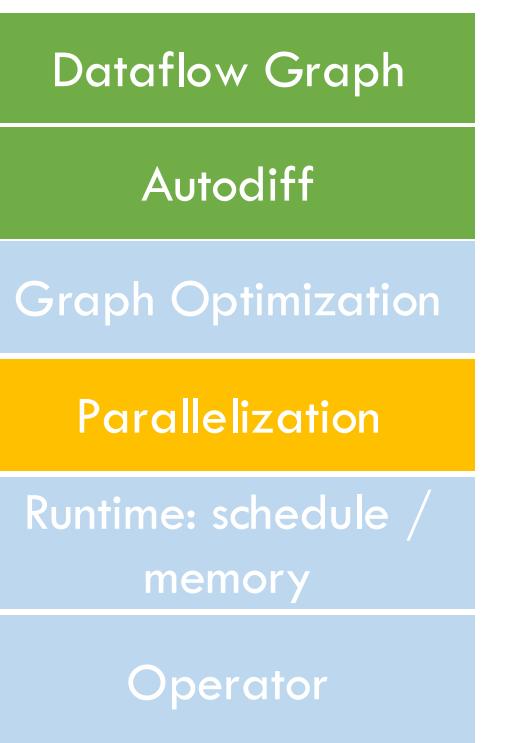


# Parallelization

- Goal: parallelize the graph compute over multiple devices

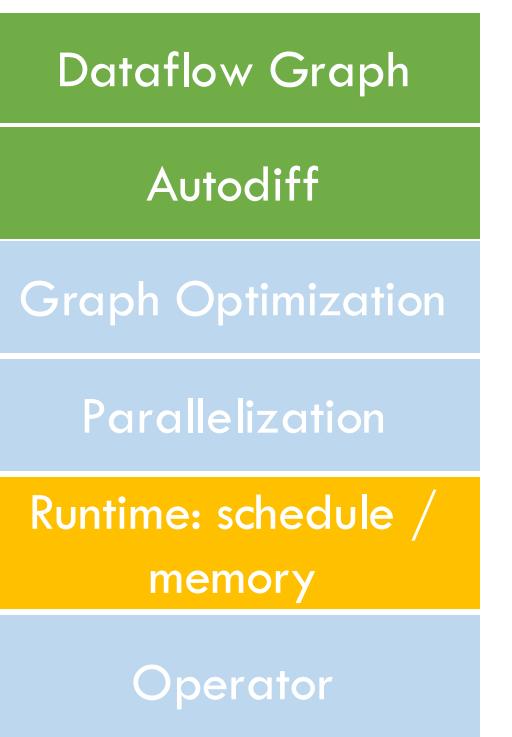
How to partition the computational graph  
on the device cluster?





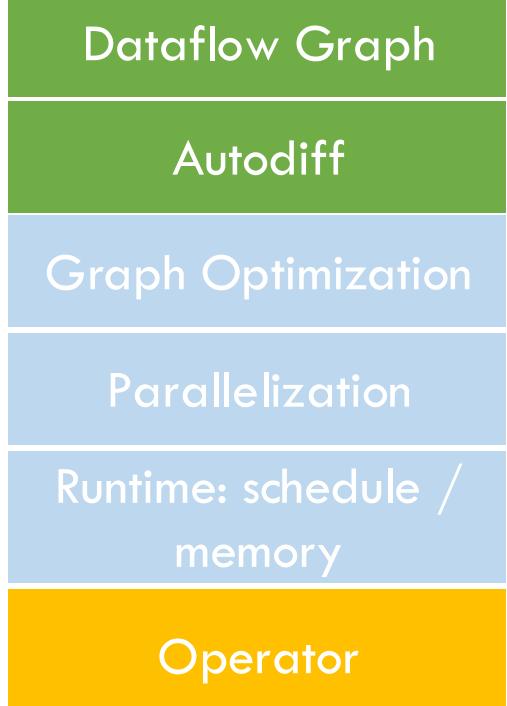
# Parallelization Problems

- How to partition
- How to communicate
- How to schedule
- Consistency
- How to auto-parallelize?



# Runtime and Scheduling

- Goal: schedule the compute/communication/memory in a way that
  - As fast as possible
  - Overlap communication with compute
  - Subject to memory constraints

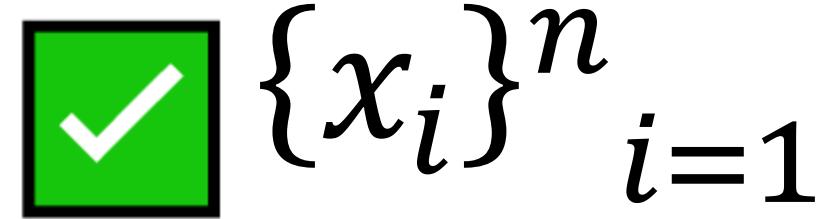


# Operator Implementation

- Goal: get the fastest possible implementation of
  - Matmul
  - Conv2d?
- For different hardware: V100, A100, H100, phone, TPU
- For different precision: fp32, fp16, fp8, fp4
- For different shape: conv2d\_3x3, conv2d\_5x5, matmul2D, 3D, attention

# High-level Picture

Data



Model

Math primitives  
(mostly matmul)



A repr that expresses the computation using primitives

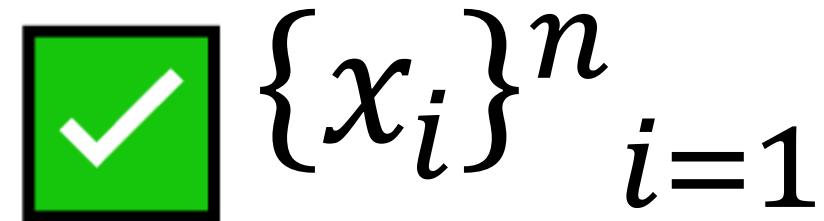
Compute

?

Make them run on (clusters of ) different kinds of hardware

# Focus of the rest of lectures

Data



Model

Math primitives  
(mostly matmul)



A repr that expresses the computation using primitives

Compute

? Make LLMs run on  
**(large clusters of ) GPUs**

# Large Language Models

- Transformers, Attentions, MoEs
- Parallelization
- Attention optimization
- Serving and inference

# Next Token Prediction

$$P(next\ word \mid prefix)$$

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

Probability("San Diego has very nice weather")  
= P("San Diego") P("has" | "San Diego")P("very" | "San Diego  
has")P("city" | ...)...P("weather" | ...)

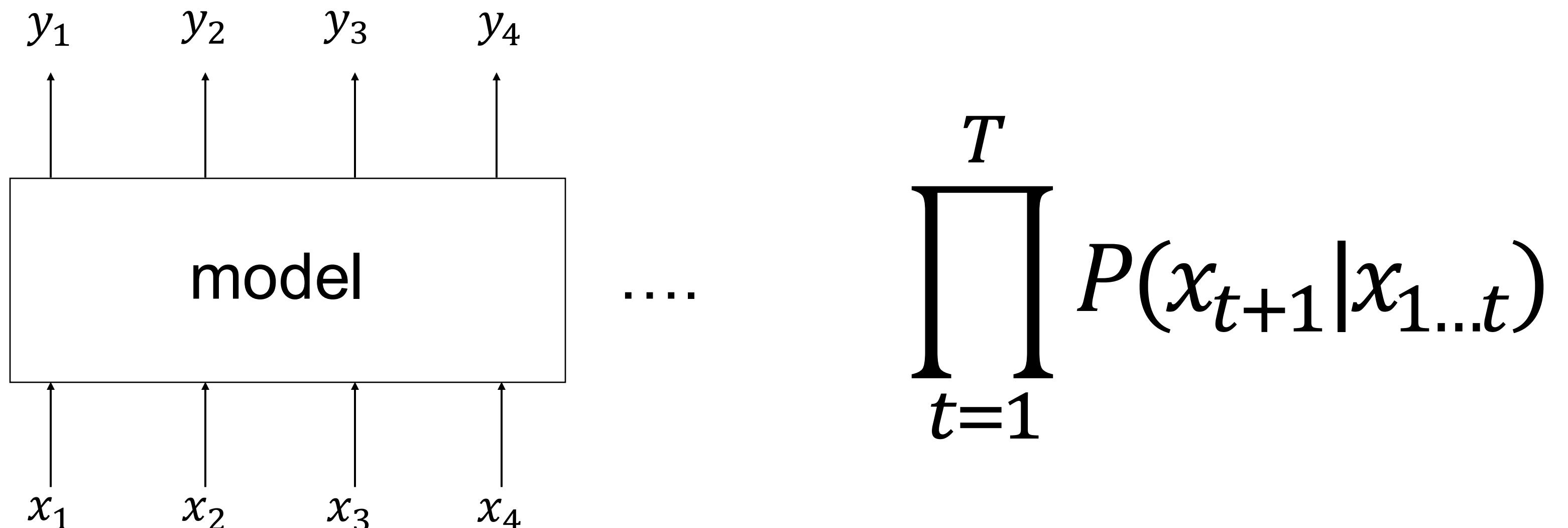
$$\text{Max Prob}(x_{1:T}) = \prod_{t=1}^T P(x_{t+1} | x_{1:t})$$

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



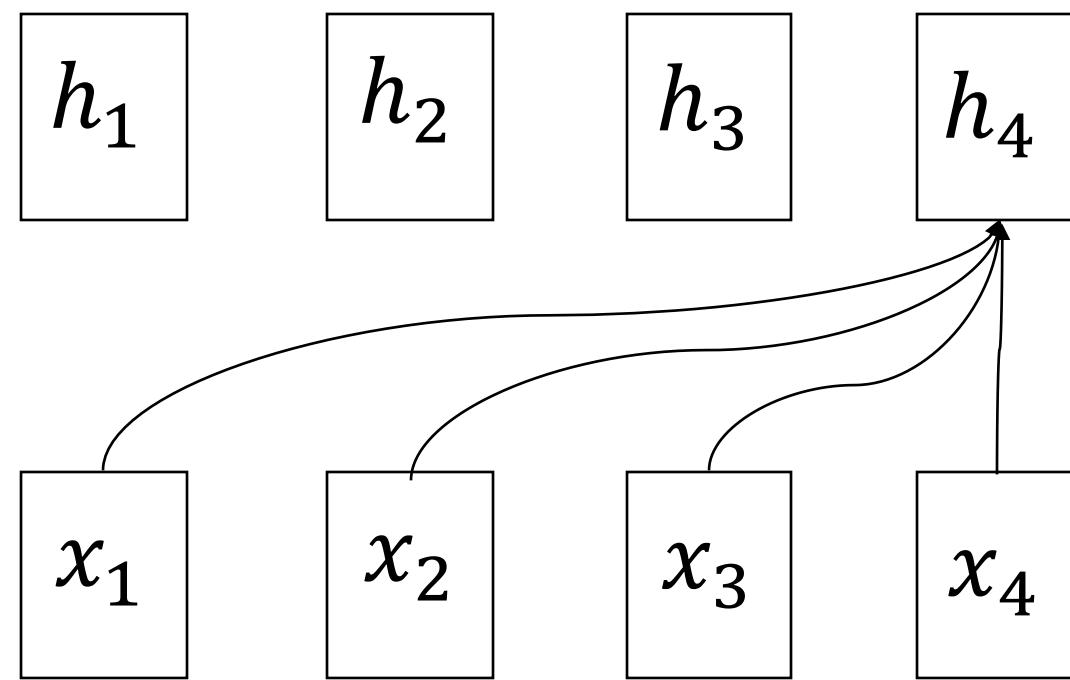
Predict each output based on history       $y_t = f_\theta(x_{1:t})$

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

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.

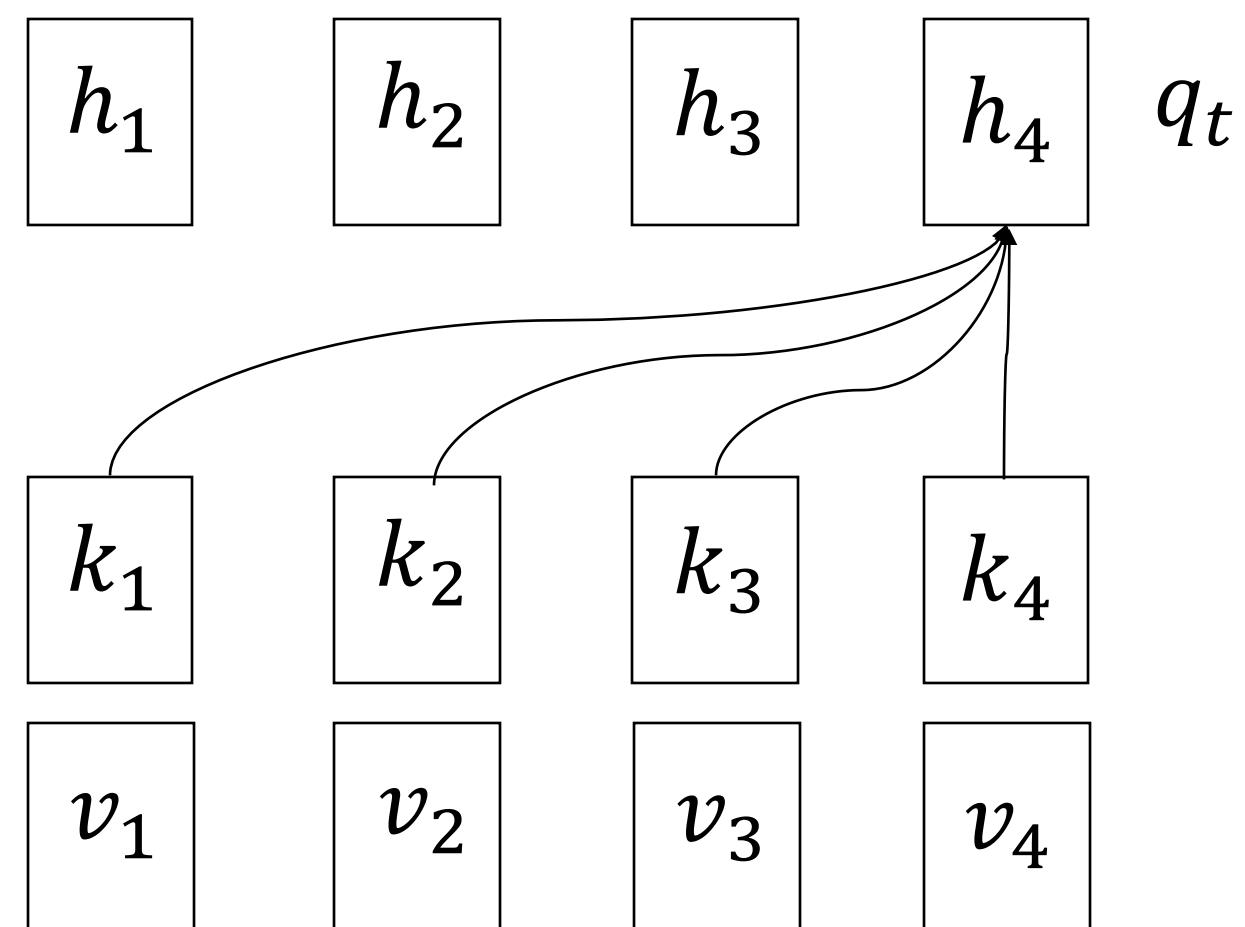
Given three inputs  $Q, K, V \in \mathbb{R}^{T \times d}$  (“queries”, “keys”, “values”)

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



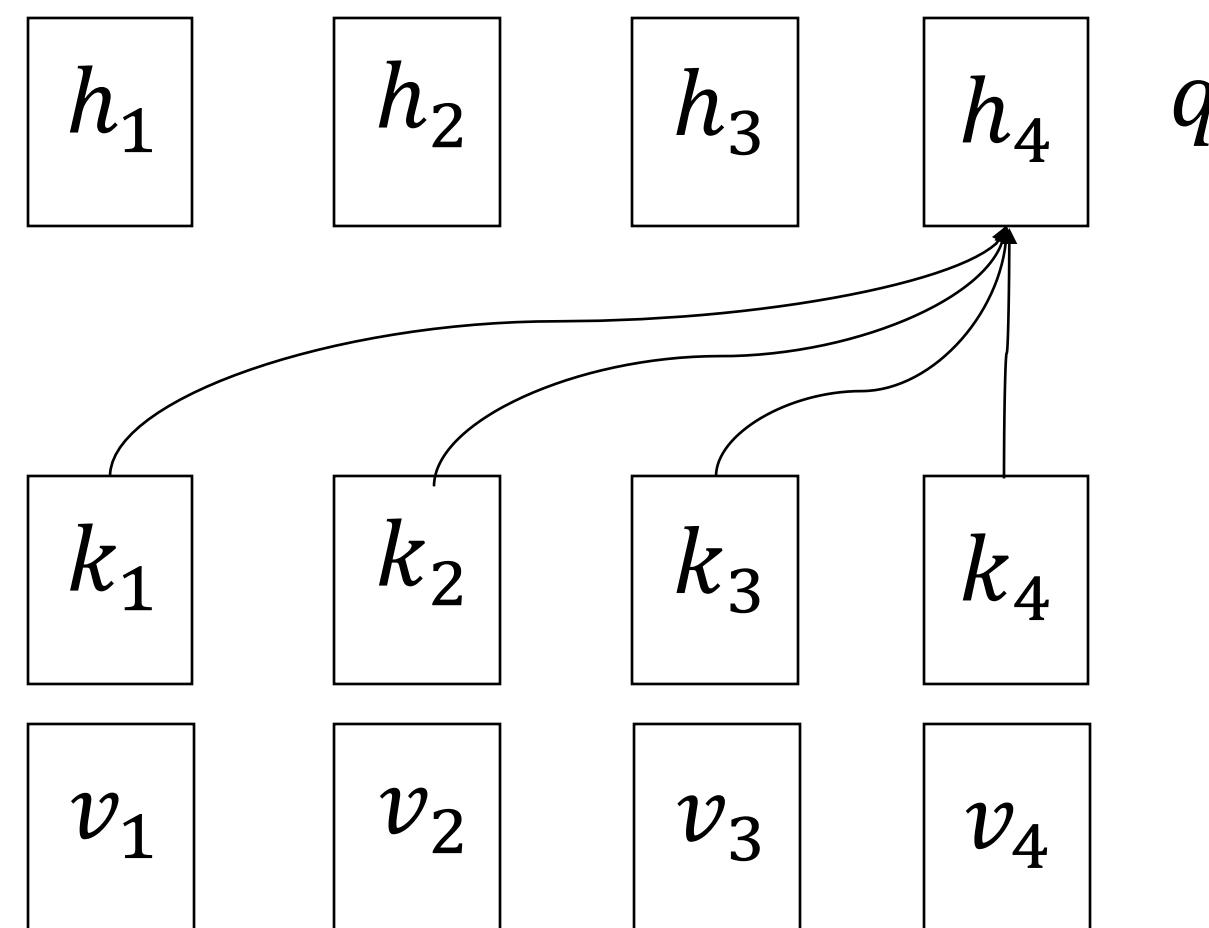
Ask the following question:

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”

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

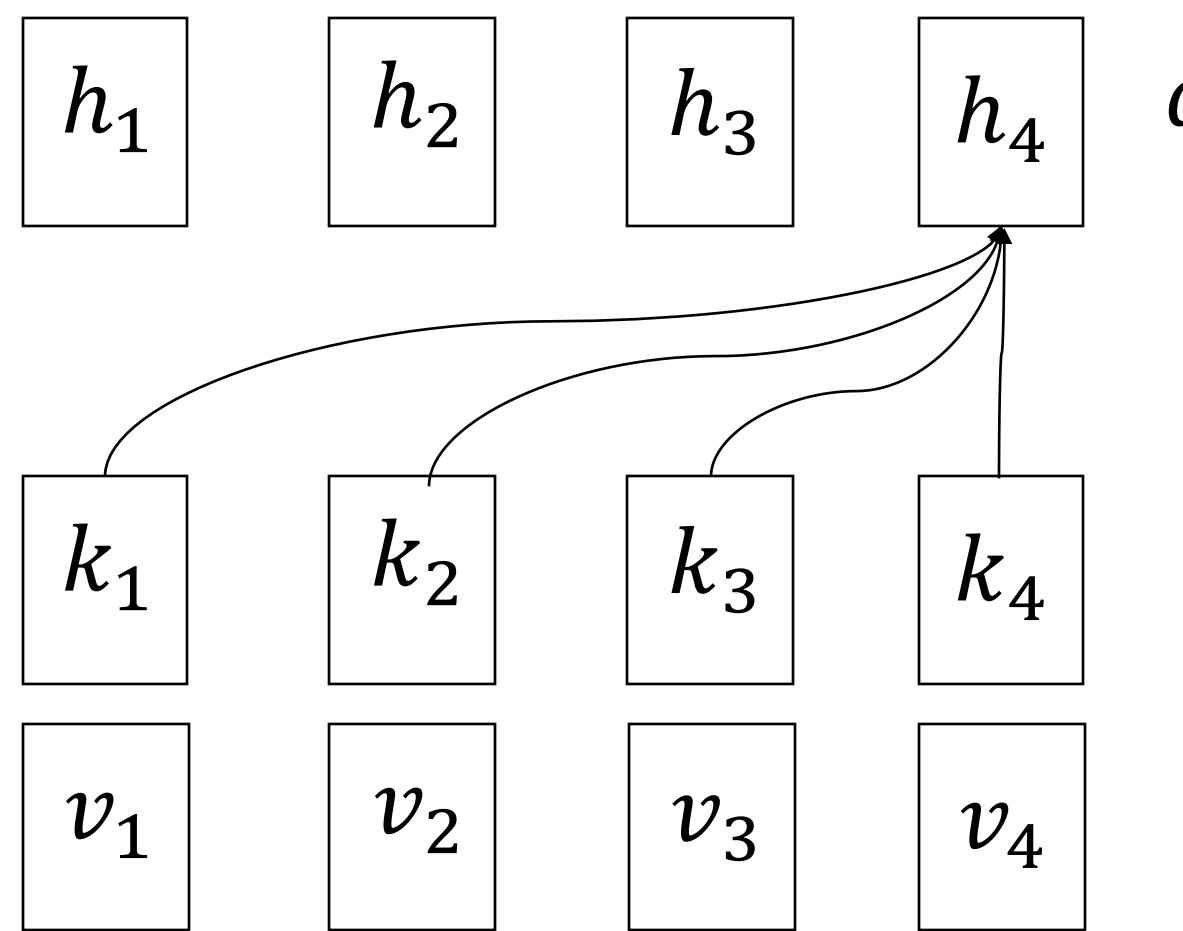
Weighed average via softmax

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

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”

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

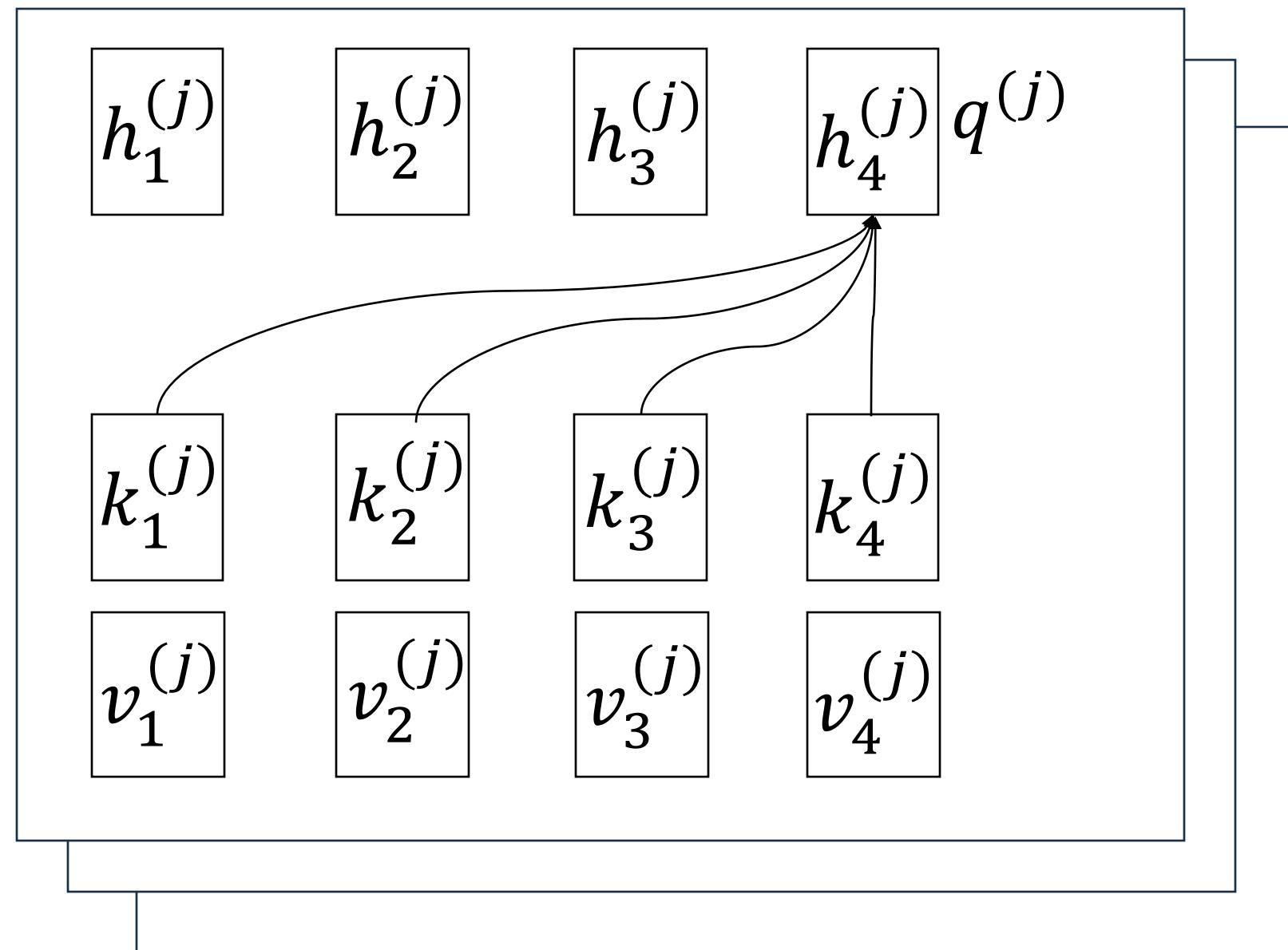
Weighed average via softmax

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

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

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

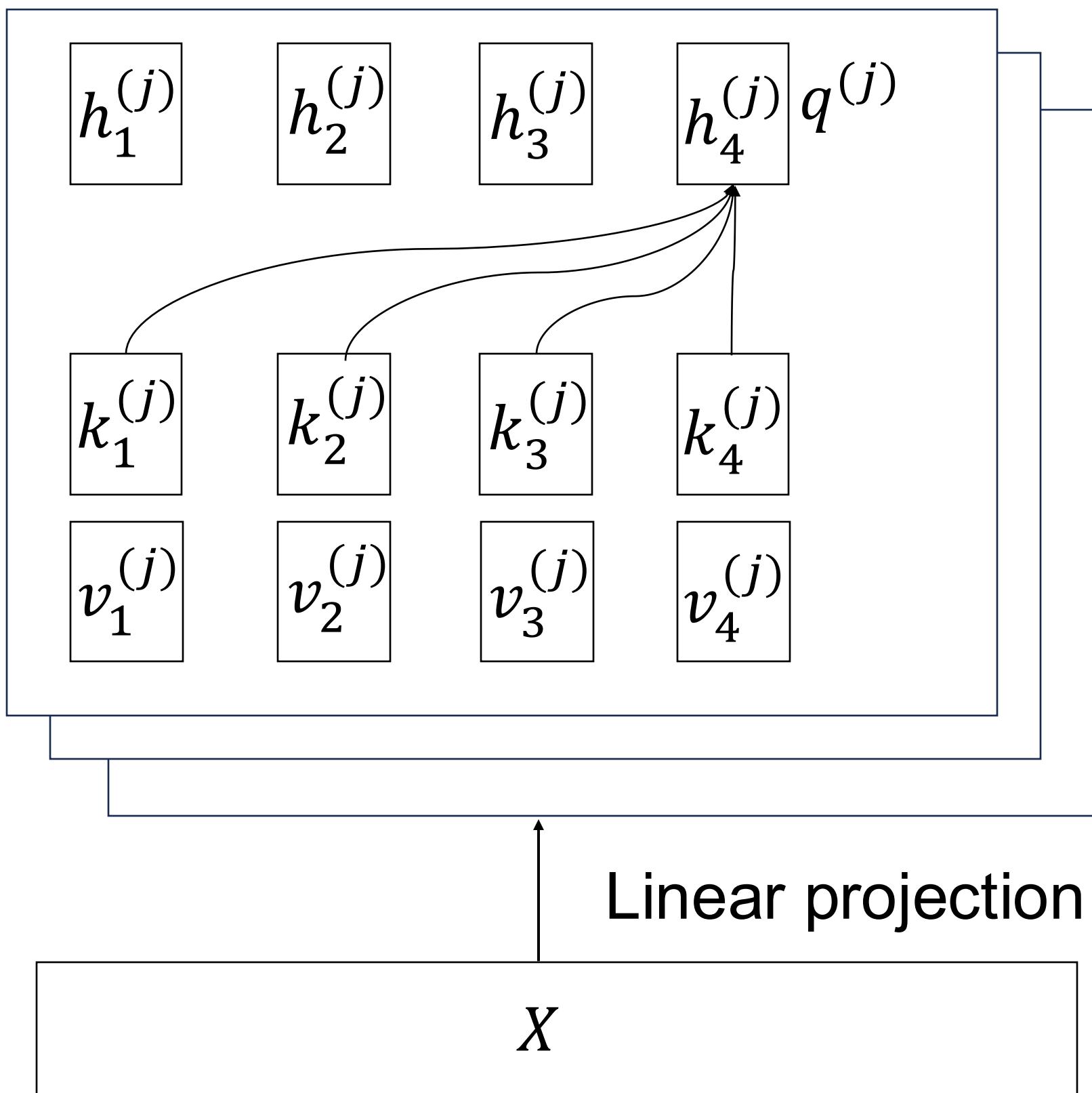
Concatenate all output heads together as output

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

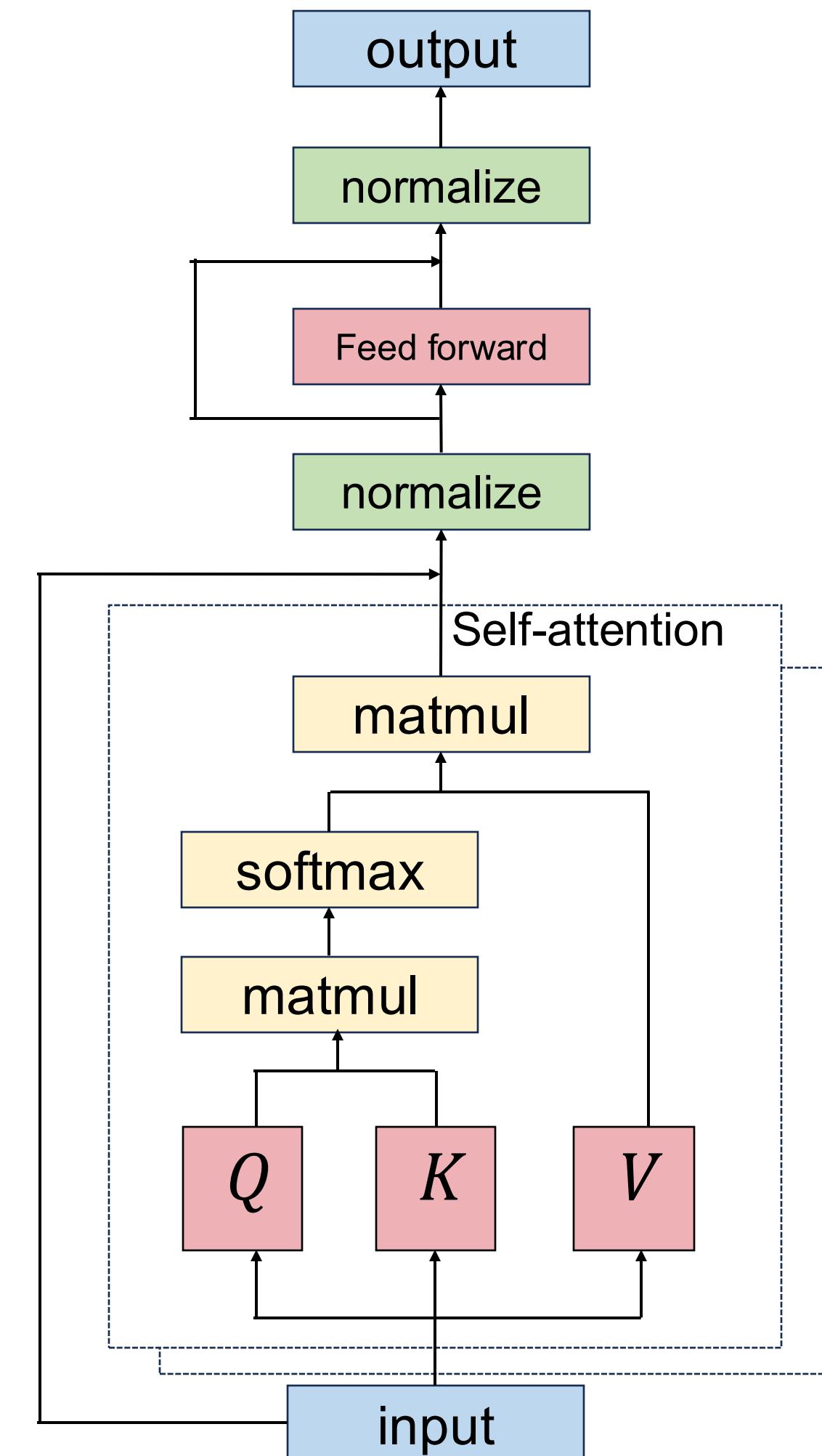
A typical transformer block

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

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

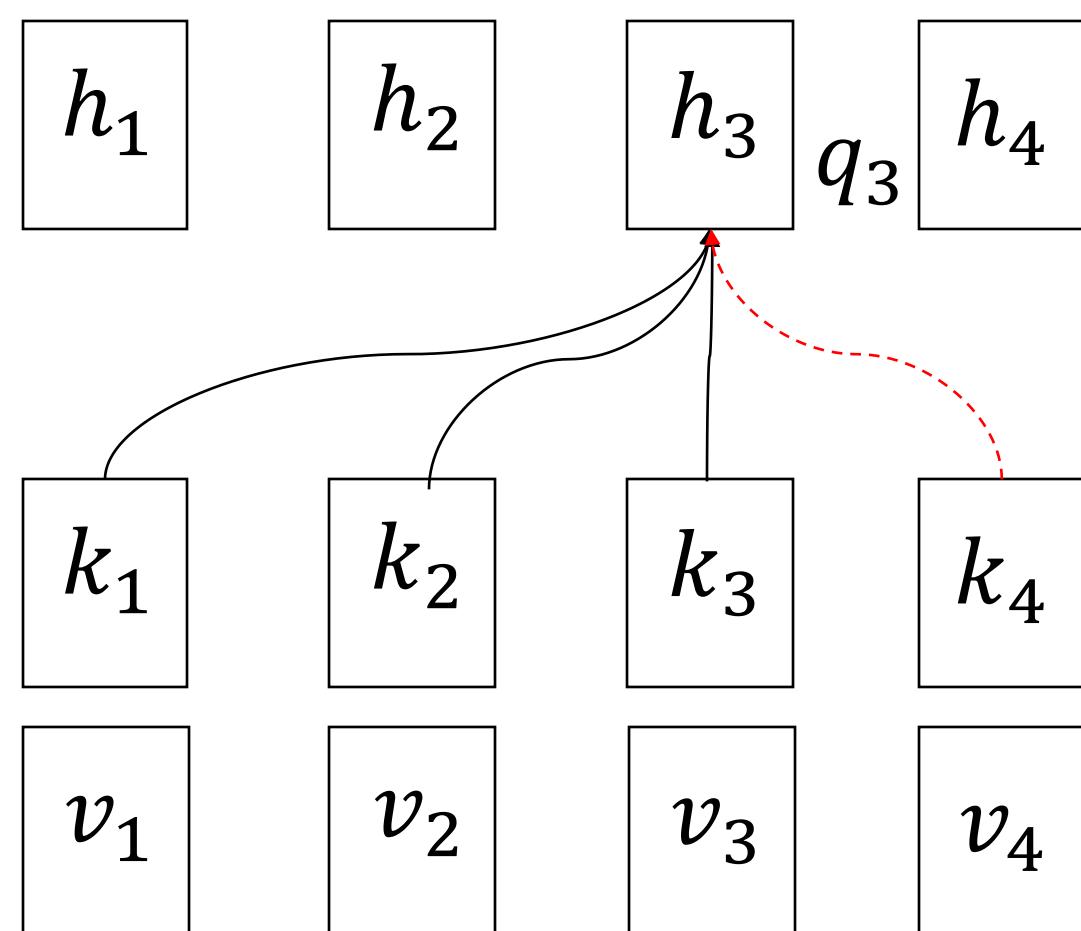
$$H = \text{LayerNorm}(\text{ReLU}(ZW_1)W_2 + Z)$$

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



# Masked Self-Attention

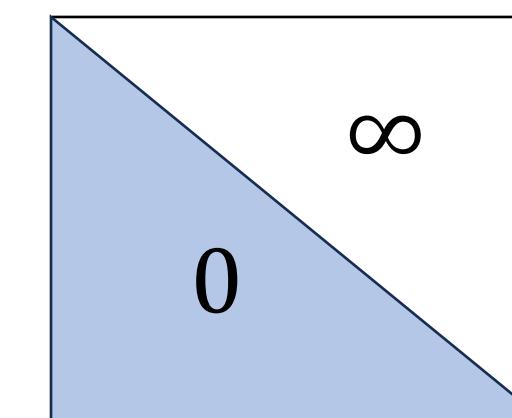
In the matrix form, we are computing weighted average over all inputs



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

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

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

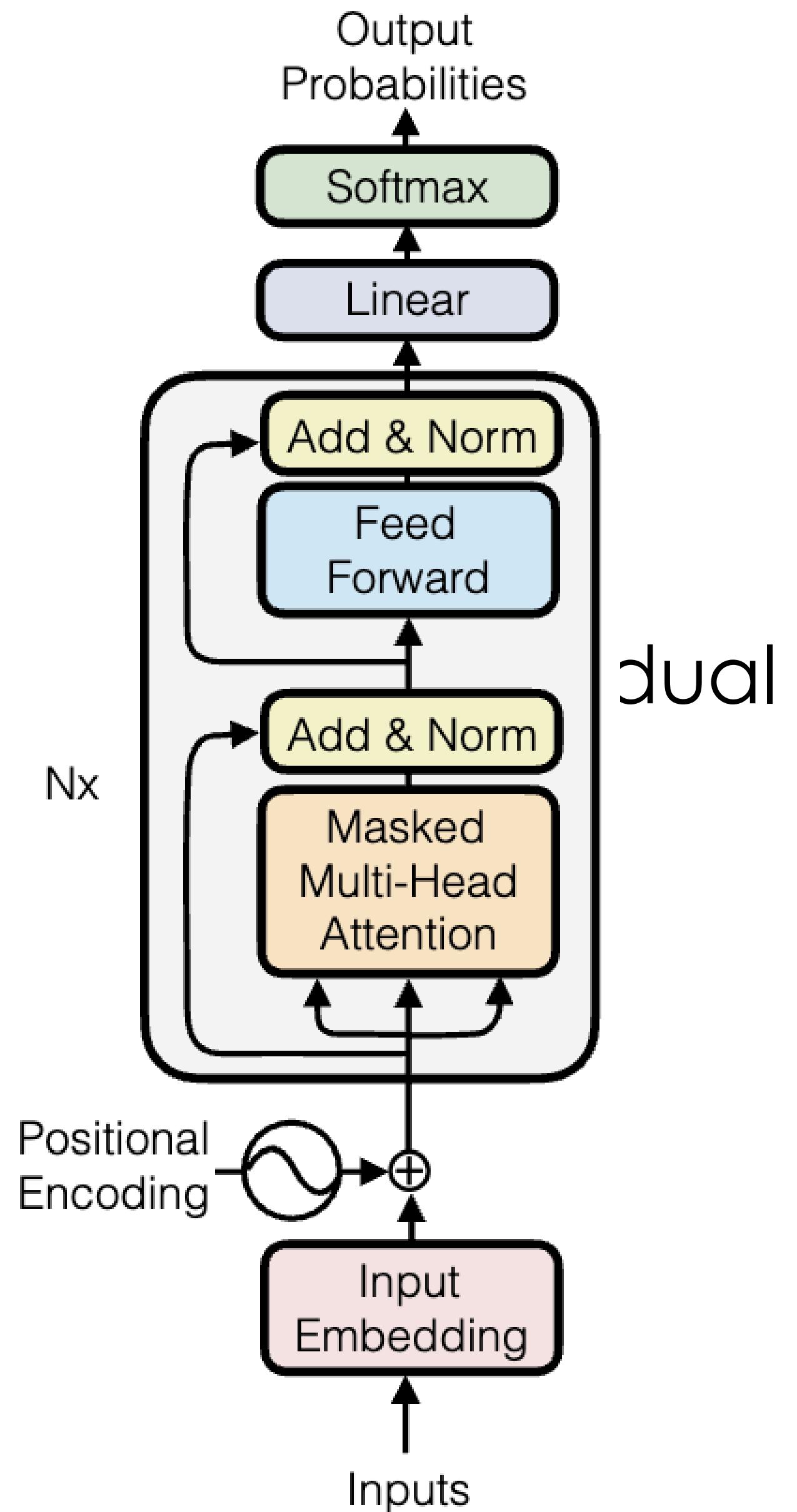


Only attend to previous inputs. Depending on input structure and model, attention mask can change.

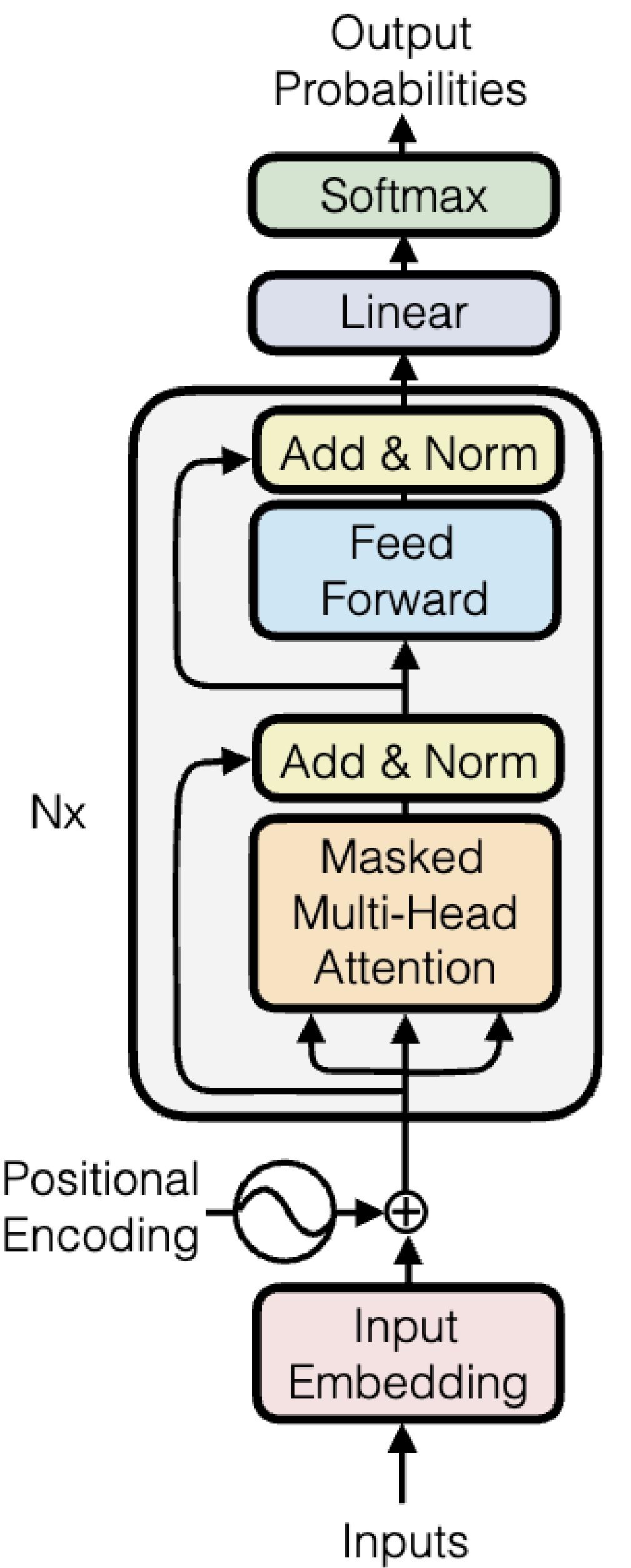
We can also simply skip the computation that are masked out if there is a special implementation to do so

# Transformers

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



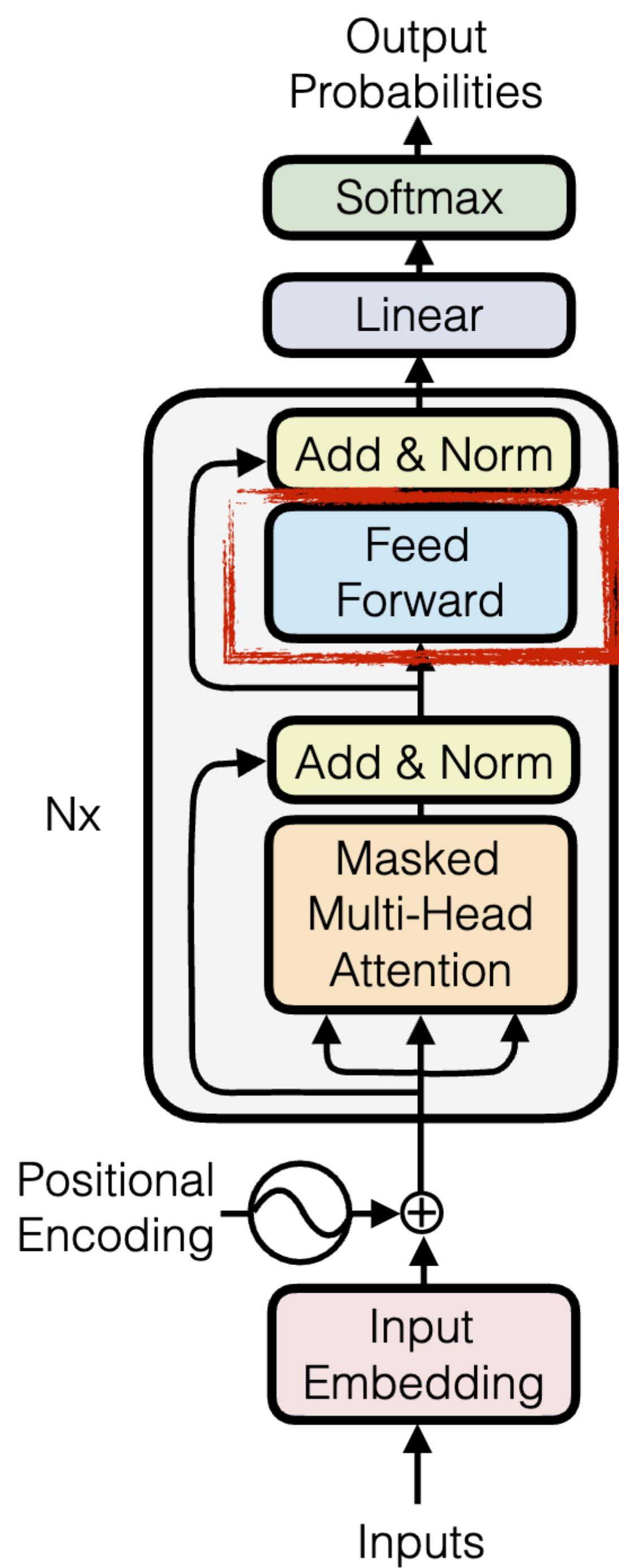
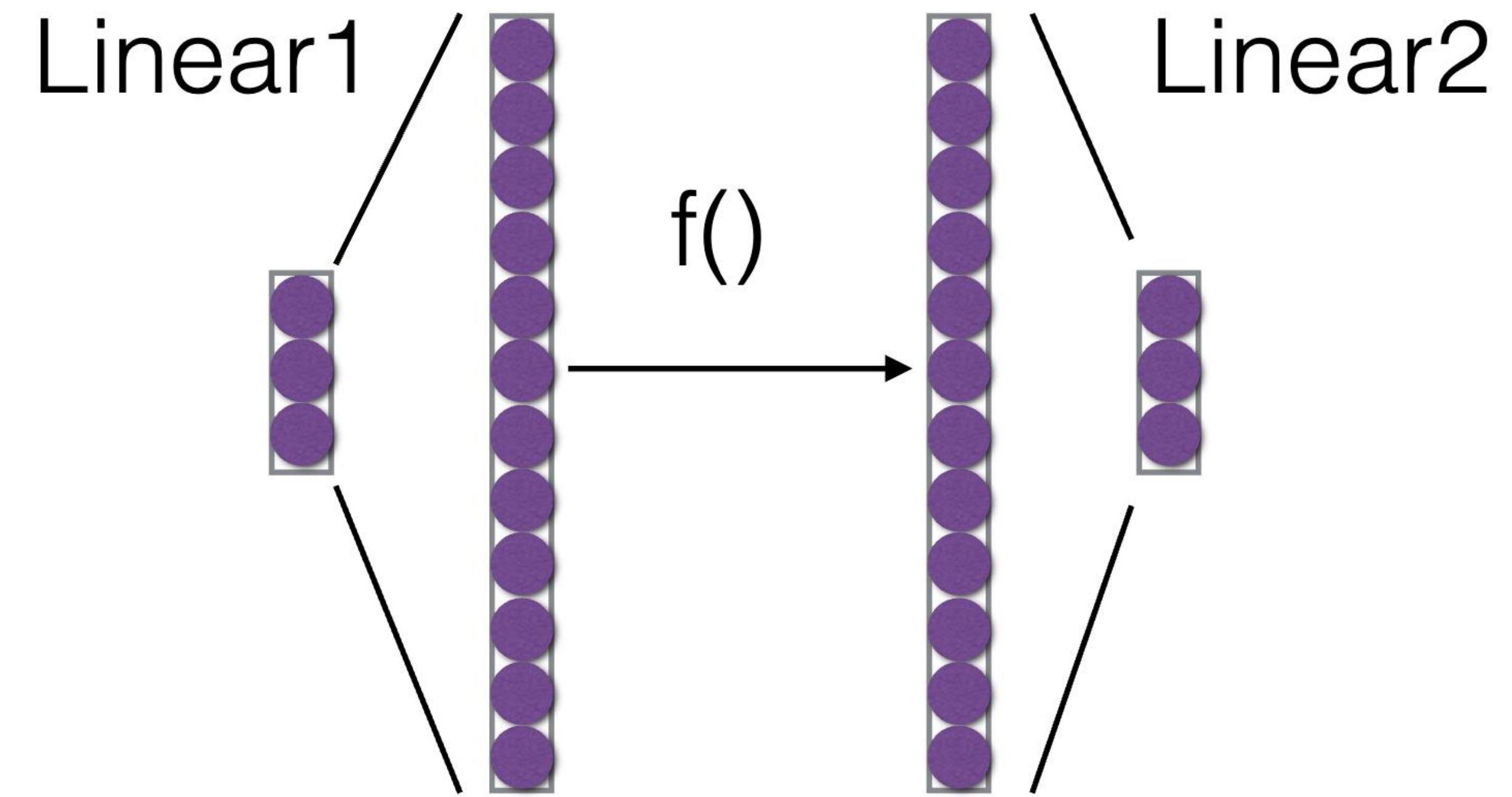
# Transformers



# Feedforward Layers

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

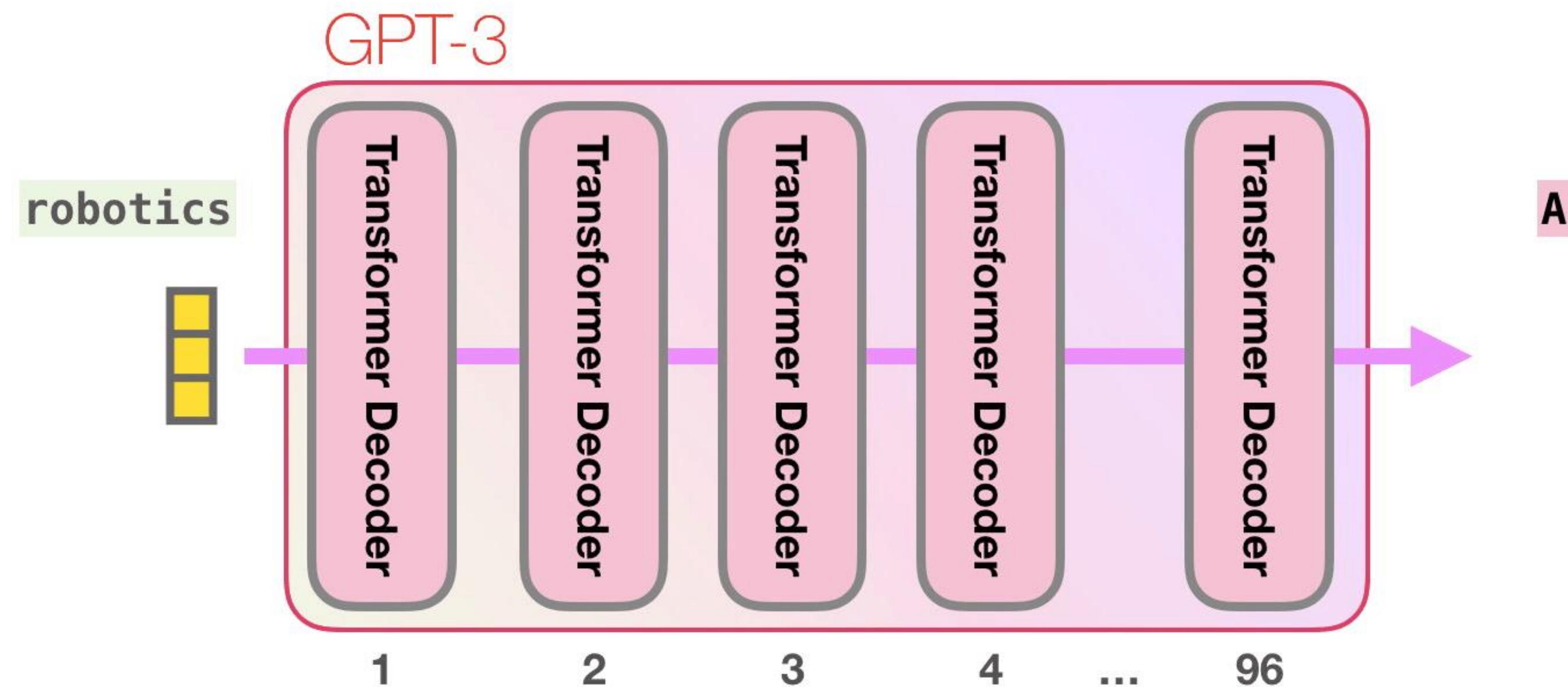
Non-linearity



# 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)
- Loss function: cross entropy loss over a sequence of words

LLMs

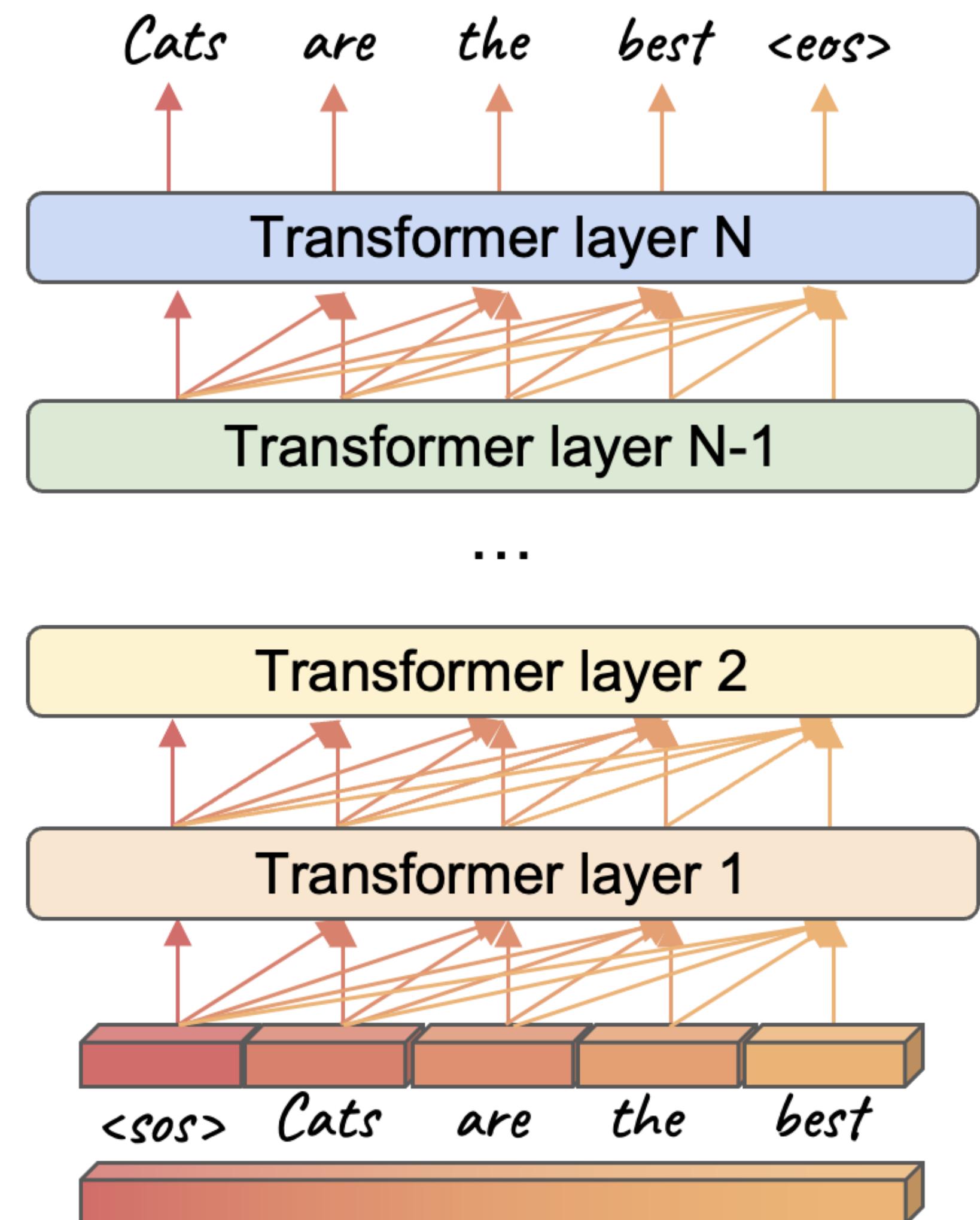


# Original Transformer vs. LLM today

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

# Training LLMs

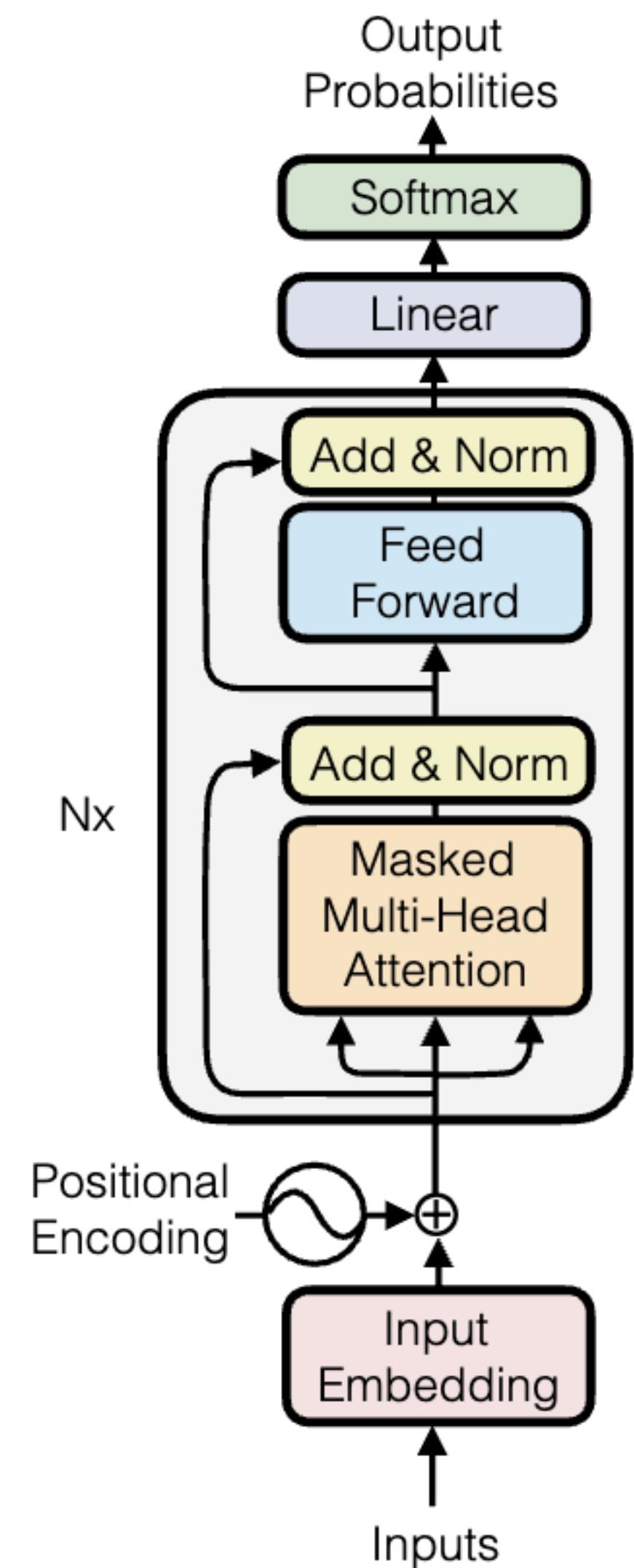
- Sequences are **known a priori**
- For each position, look at  $[1, 2, \dots, t-1]$  words to predict word  $t$ , and calculate the loss at  $t$
- Parallelize the computation across all token positions, and then apply masking



# Connecting the Dots: Compute/Comm characteristic of LLMs

Key characteristics: compute, memory, communication

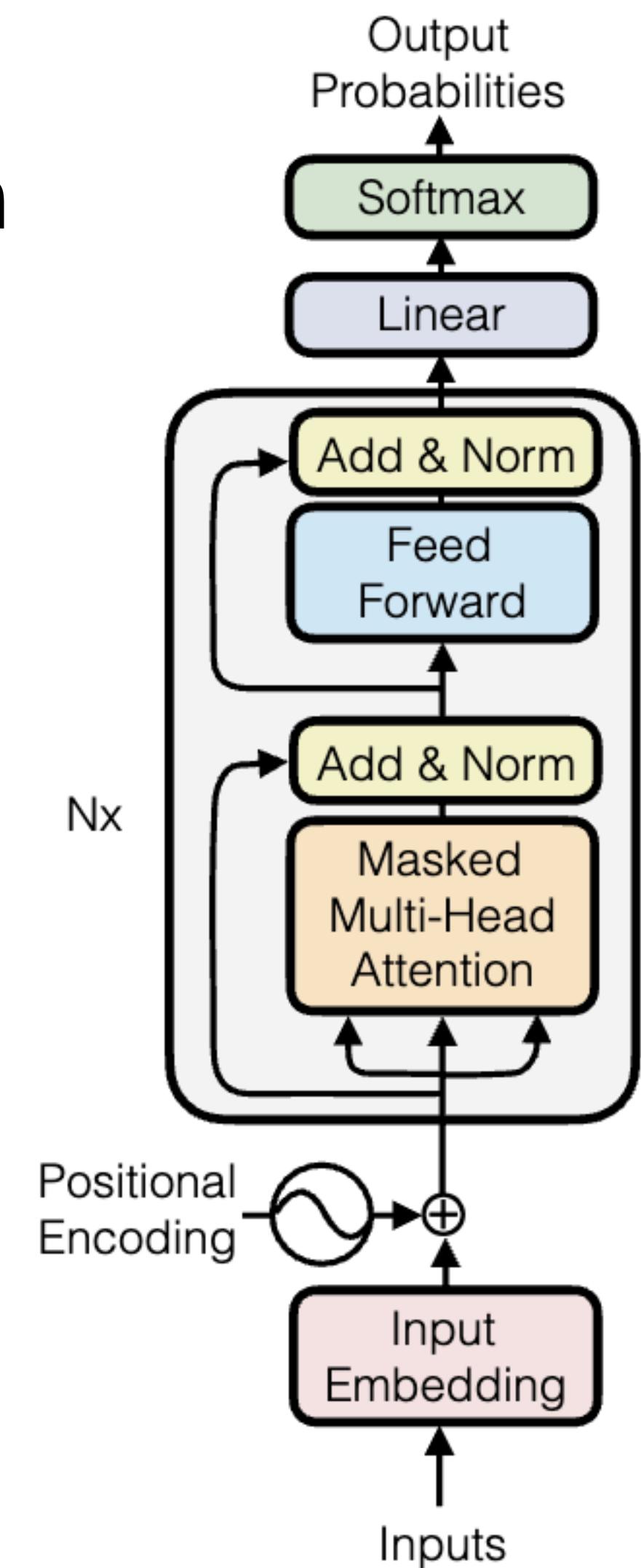
- calculate the number of parameters of an LLM?
  - memory, communication
- calculate the flops needed to train an LLM?
  - compute
- calculate the memory needed to train an LLM?
  - memory, communication

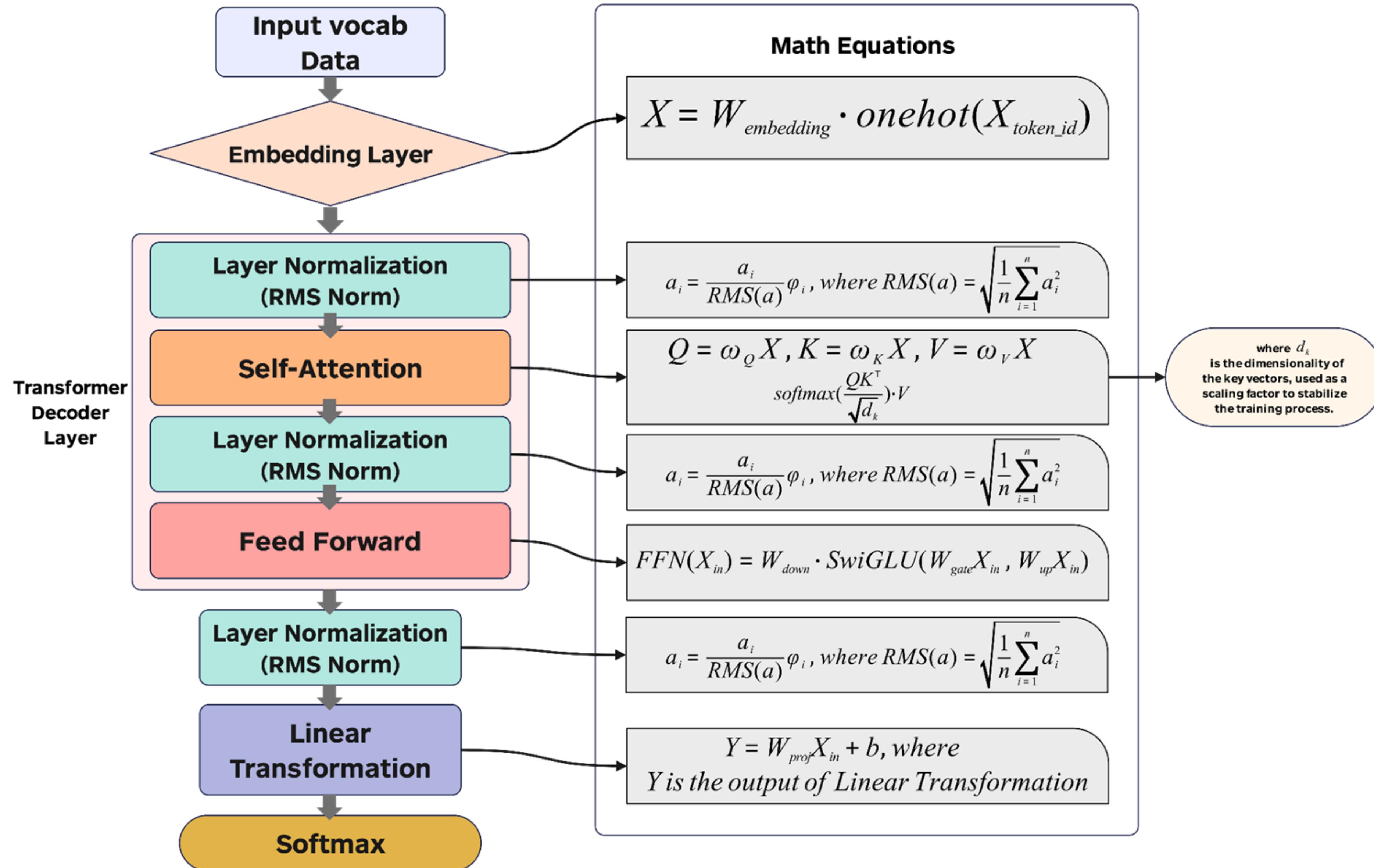


# Connecting the Dots: Compute/Comm characteristic of LLMs

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

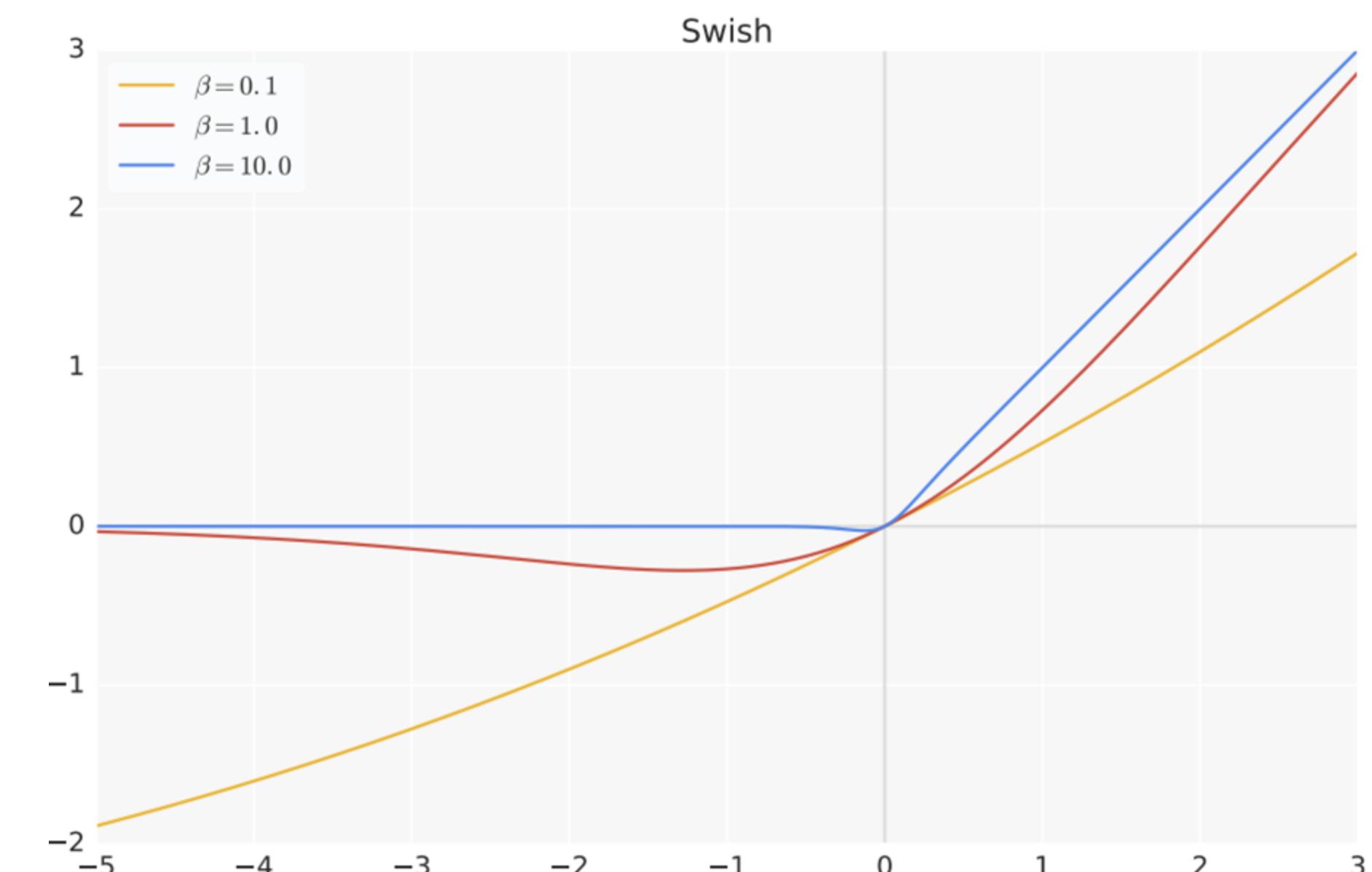
The general formula for SwiGLU is:

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

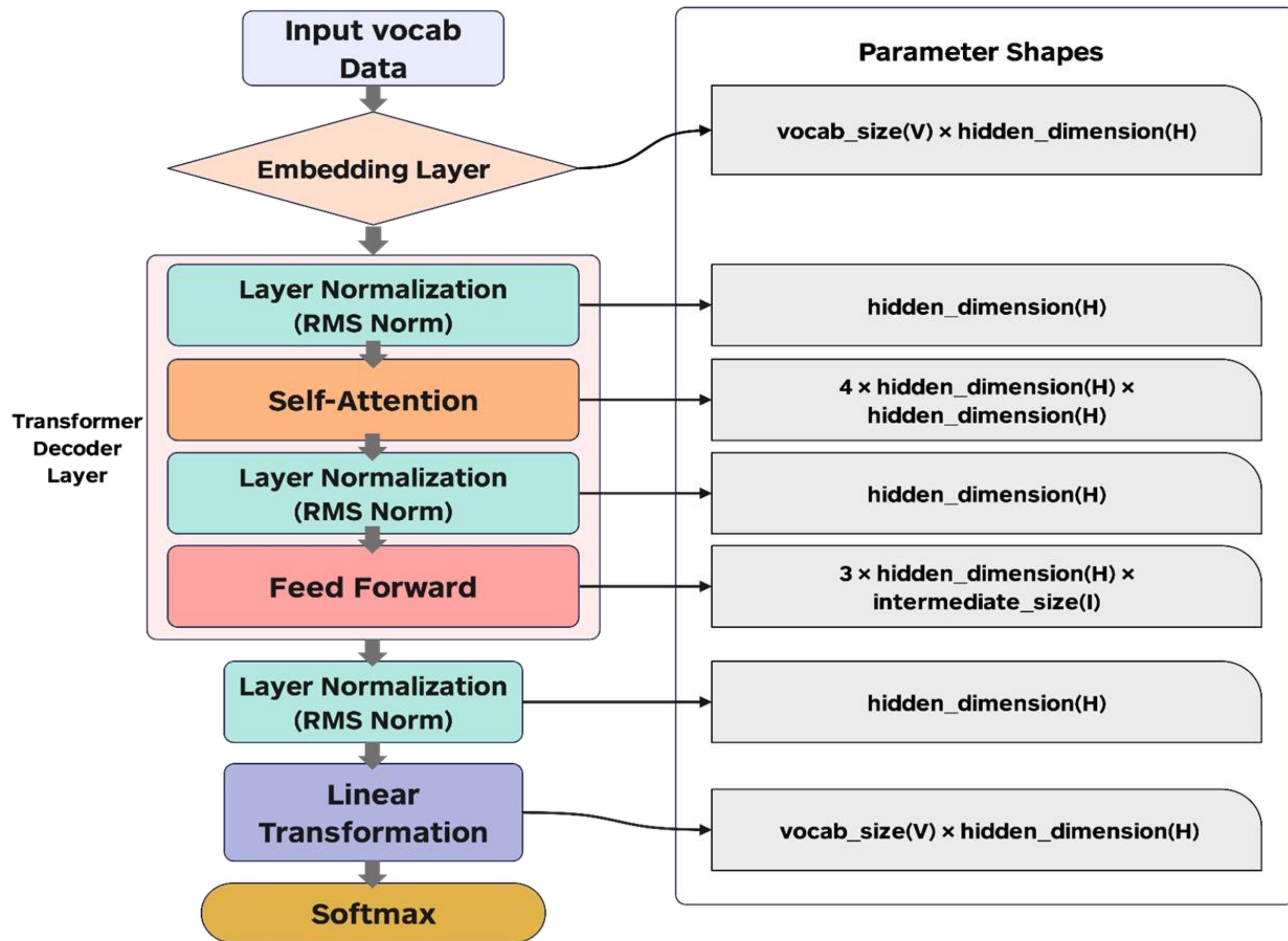
**Swish** is the activation function applied to one branch, defined as:

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

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



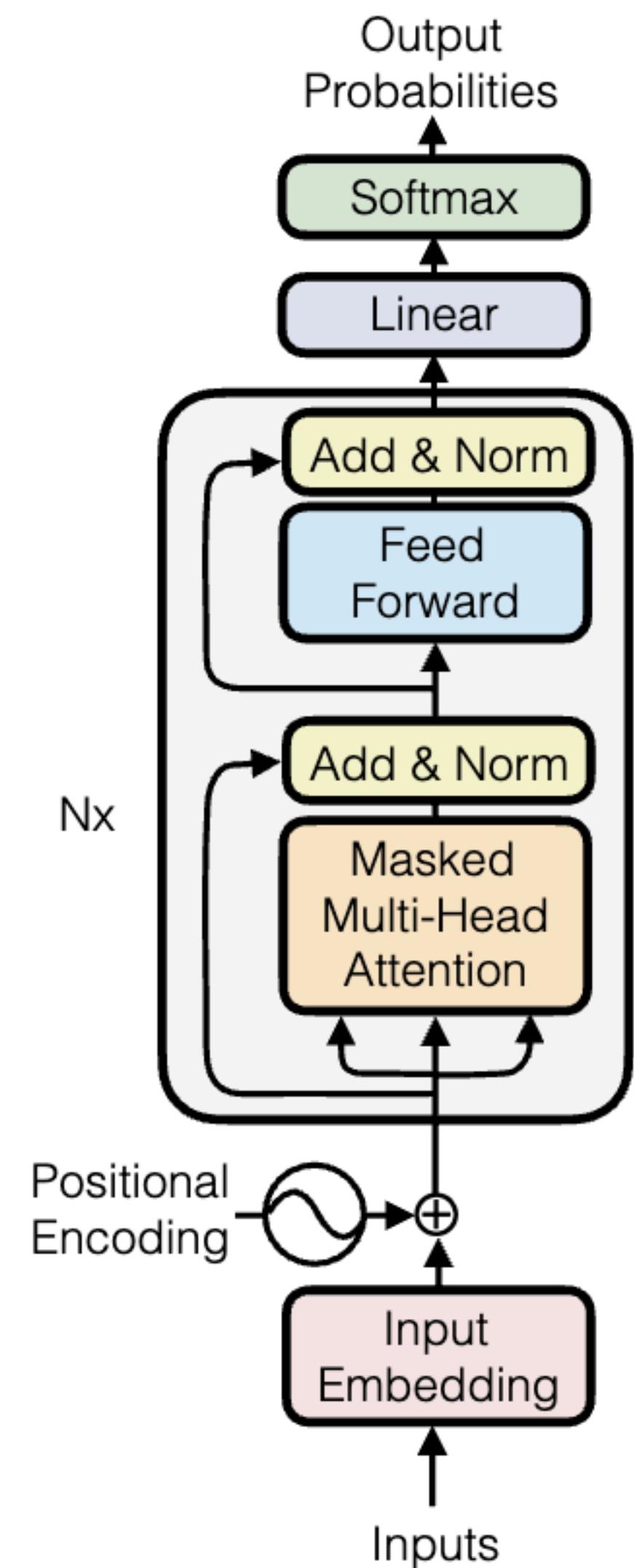
# Summary



# 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?

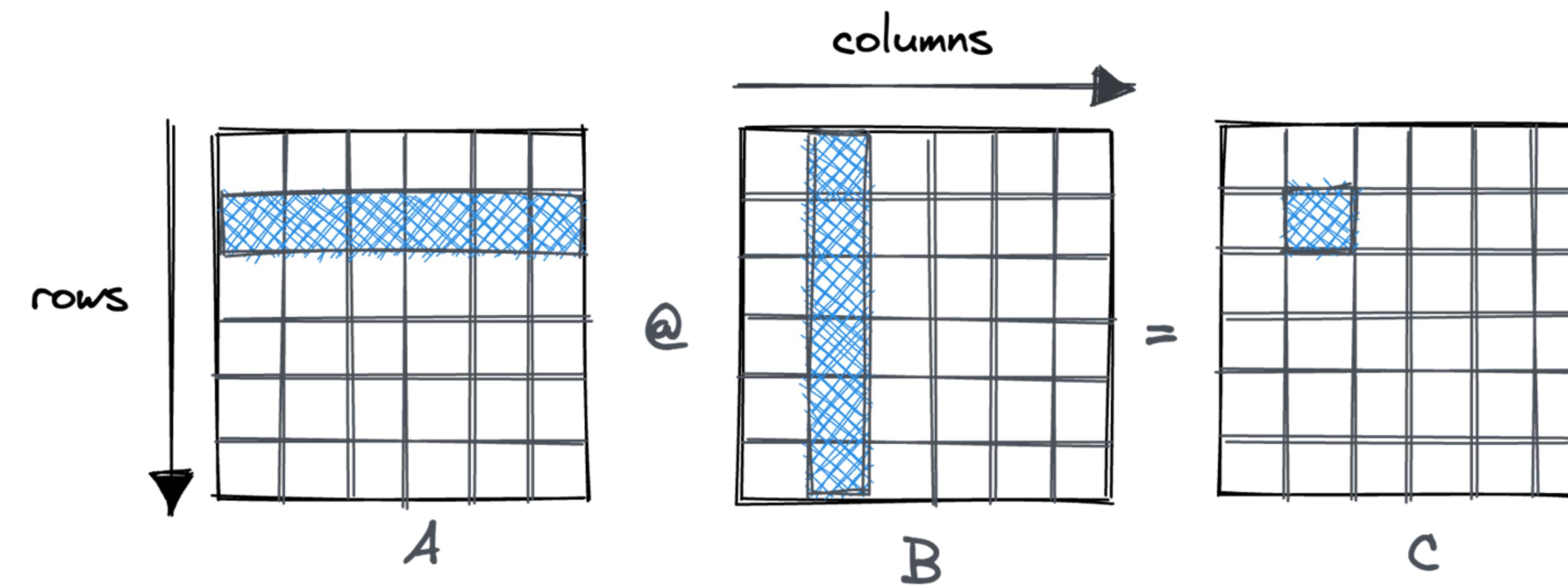


# Estimate the Compute: FLOPs

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:

$X$

Self Attention:

$XW_Q, XW_K, XW_V$

RoPE

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

PV

$AW_O$

Residual Connection:

Output Shape:

( $b, s, h$ )

( $b, s, h$ )

( $b, n, s, d$ )

( $b, n, s, s$ )

( $b, n, s, d$ )

( $b, s, h$ )

( $b, s, h$ )

FLOPs

0

$3 * 2bsh^2$

$3bsnd$

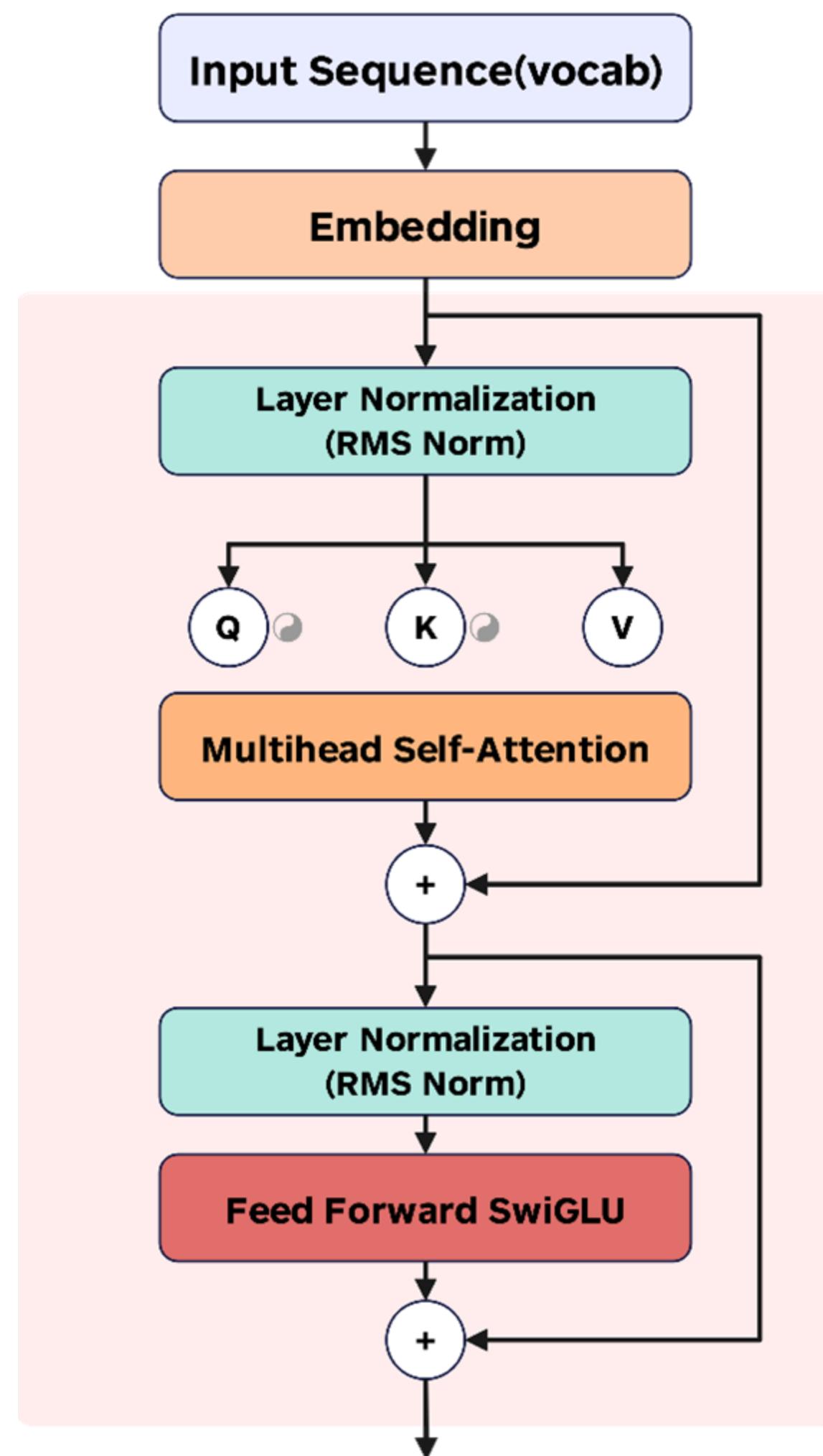
$2bs^2nd + 3bs^2n$

$2bs^2nd$

$2bsh^2$

$bsh$

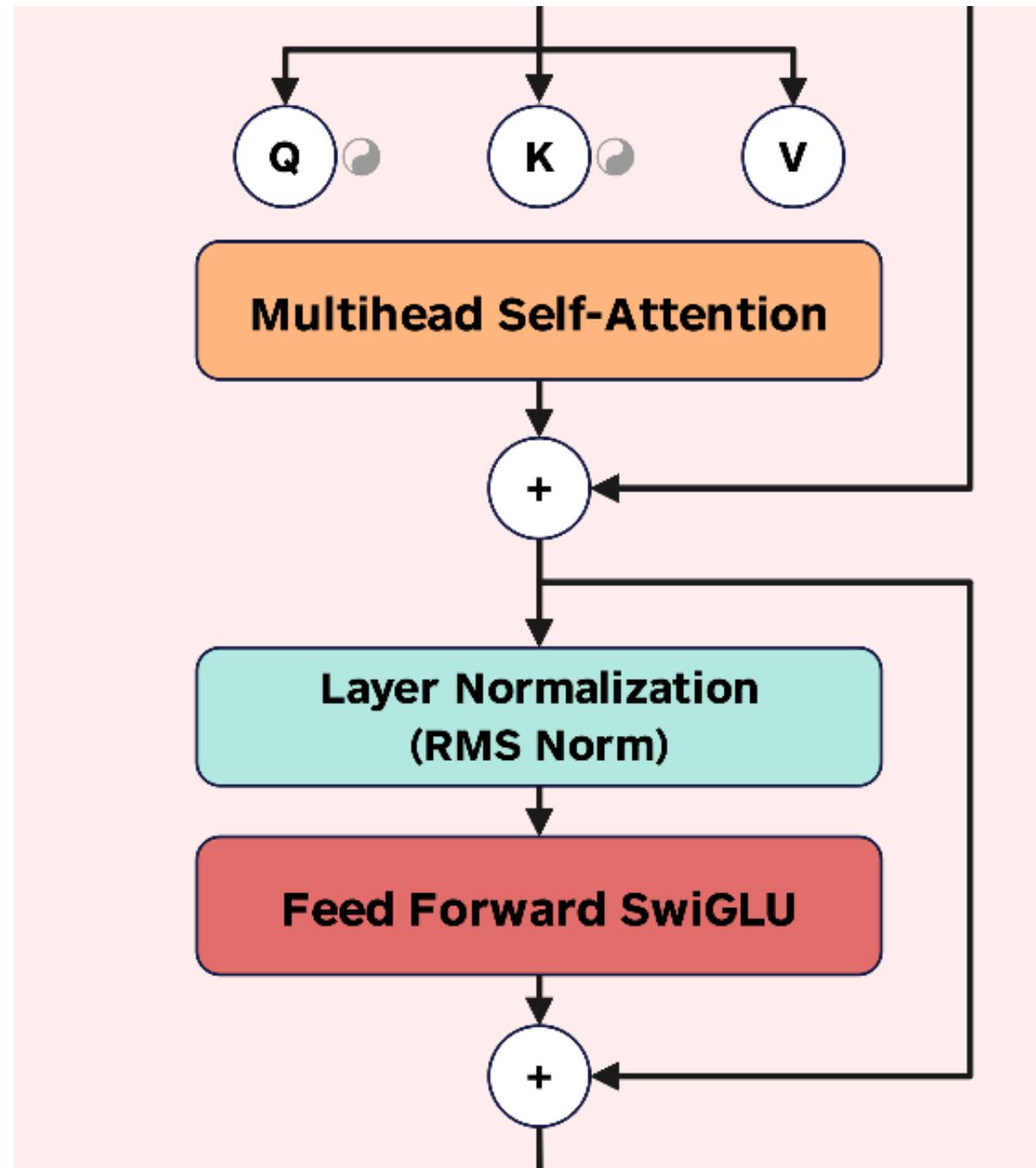
Batch size:  $b$   
Sequence length:  $s$   
# of attention heads:  $n$   
Hidden state dim of one head:  $d$   
Hidden state dim:  $h$



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

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:		
	( $b, s, h$ )	$4bsh + 2bs$

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



### 1. Calculate Root Mean Square:

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

### 2. Normalize:

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

# LLama 2 7B Flops Forward (Training)

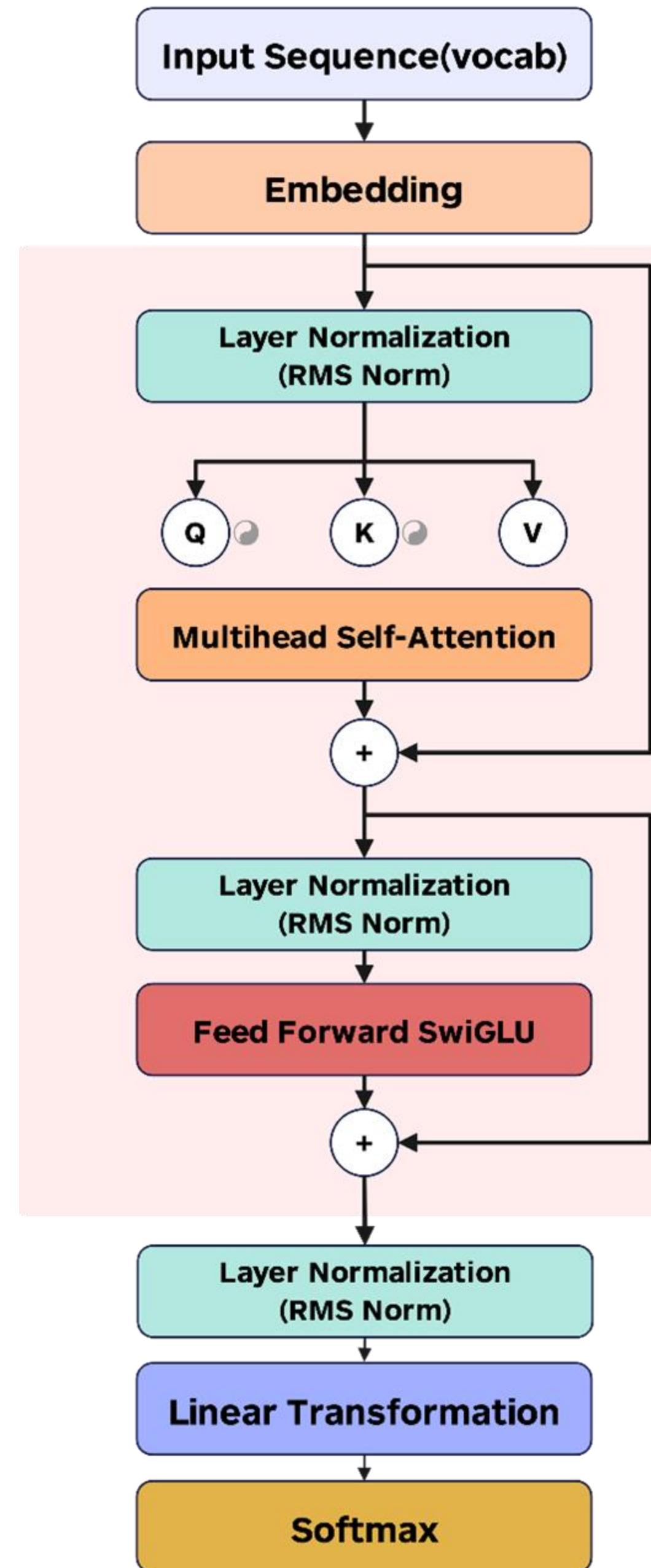
Total Flops  $\approx$  #num\_layers \* (Attention block + SwiGLU block)

+ Prediction head

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

$$+ \text{#num\_layers} ( 6bshi)$$

$$+ 2 bshv$$



# 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

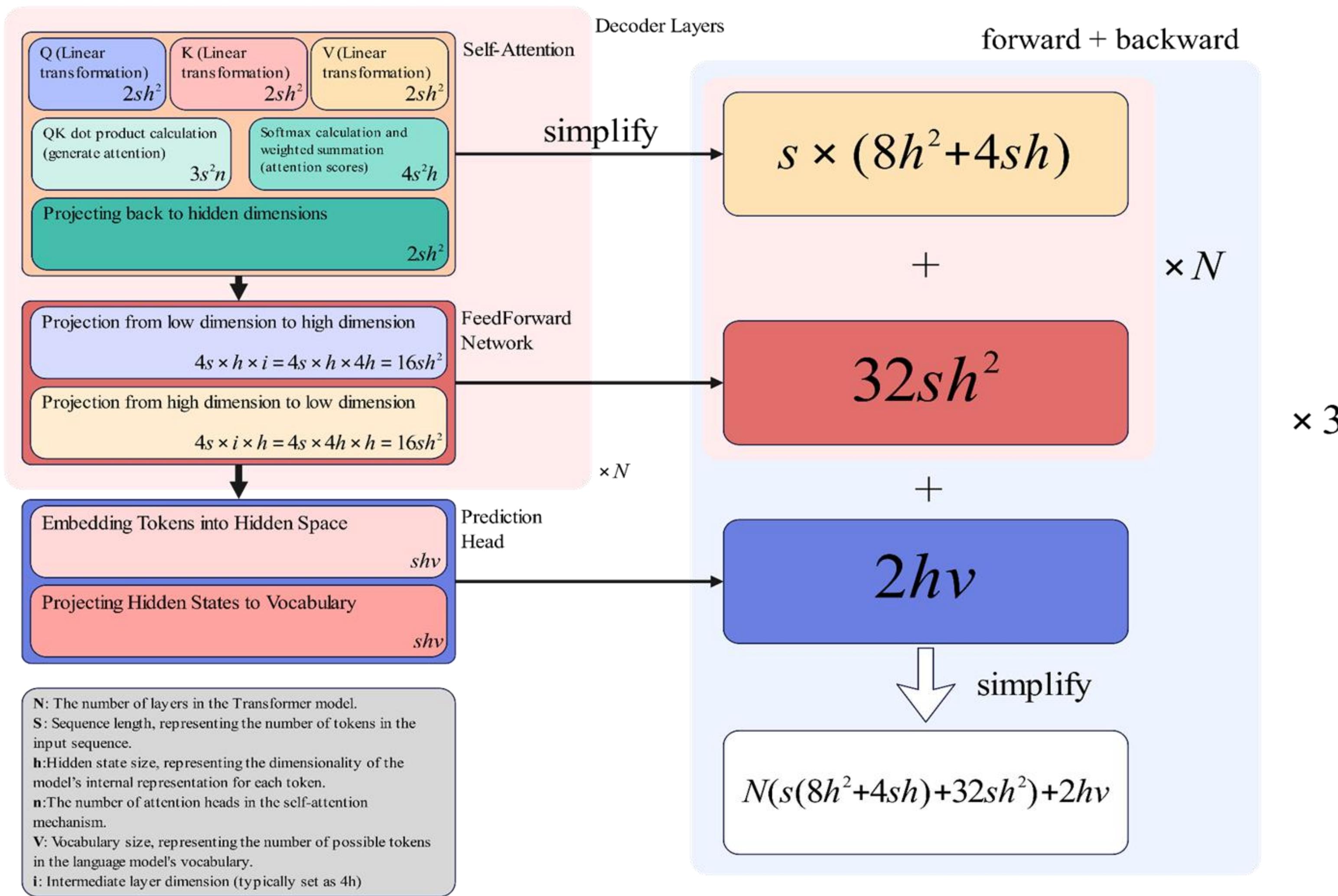
$$\begin{aligned}\text{Total Flops} &\approx N * (6bsh^2 + 4bs^2h + 3bs^2n + 2bsh^2) \\ &\quad + N (6bshi) \\ &\quad + 2 bshv \\ &\approx 63 \text{ TFLOPs}\end{aligned}$$

# Flops Distribution

## **Training Computational Costs Breakdown:**

- . **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?

