



# **LLaMA From Scratch (LLaMA 1 & LLaMA 2)**

# Prerequisites

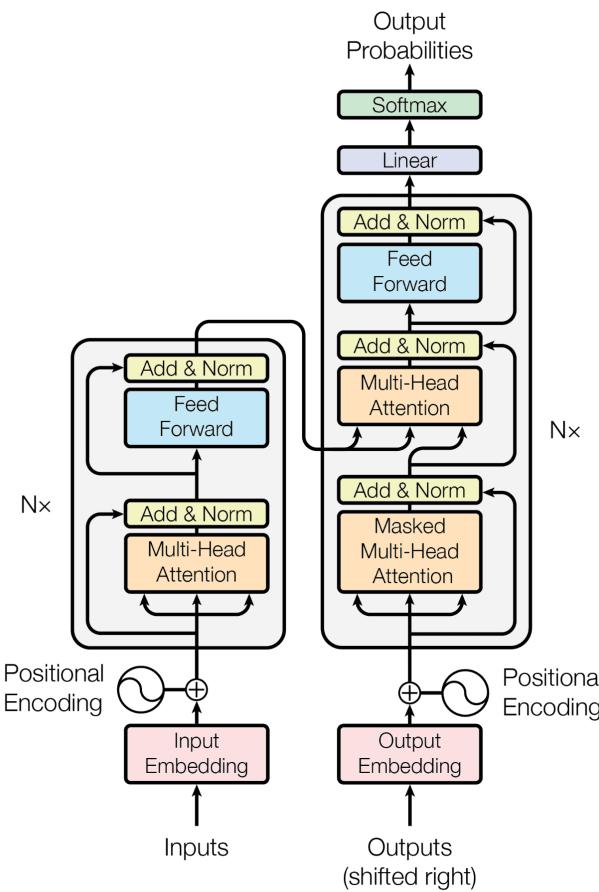
- Structure of the Transformer model and how the attention mechanism works
- Training and inference of a Transformer model
- Linear algebra
  - Matrix Multiplication
  - Dot product
- Complex numbers
  - Euler's formula (not fundamental, nice to know)

$$e^{ix} = \cos x + i \sin x$$

# Topics

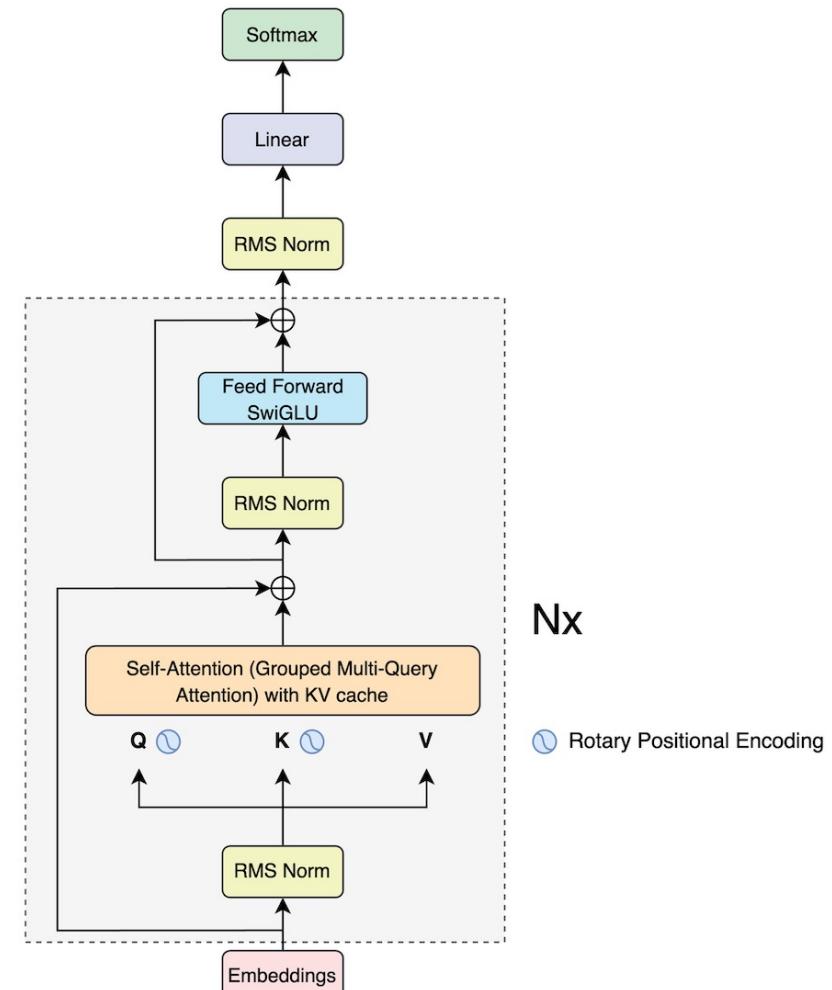
- Architectural differences between vanilla Transformer and LLaMA
- RMS Normalization (with review of Layer Normalization)
- Rotary Positional Embeddings
- KV-Cache
- Multi-Query Attention
- Grouped Multi-Query Attention
- SwiGLU Activation Function

# Transformer vs LLaMA



**Transformer**

(“Attention is All You Need”)



**LLaMA**

# Models (LLaMA 1)

params	dimension	n heads	n layers	learning rate	batch size	n tokens
6.7B	4096	32	32	$3.0e^{-4}$	4M	1.0T
13.0B	5120	40	40	$3.0e^{-4}$	4M	1.0T
32.5B	6656	52	60	$1.5e^{-4}$	4M	1.4T
65.2B	8192	64	80	$1.5e^{-4}$	4M	1.4T

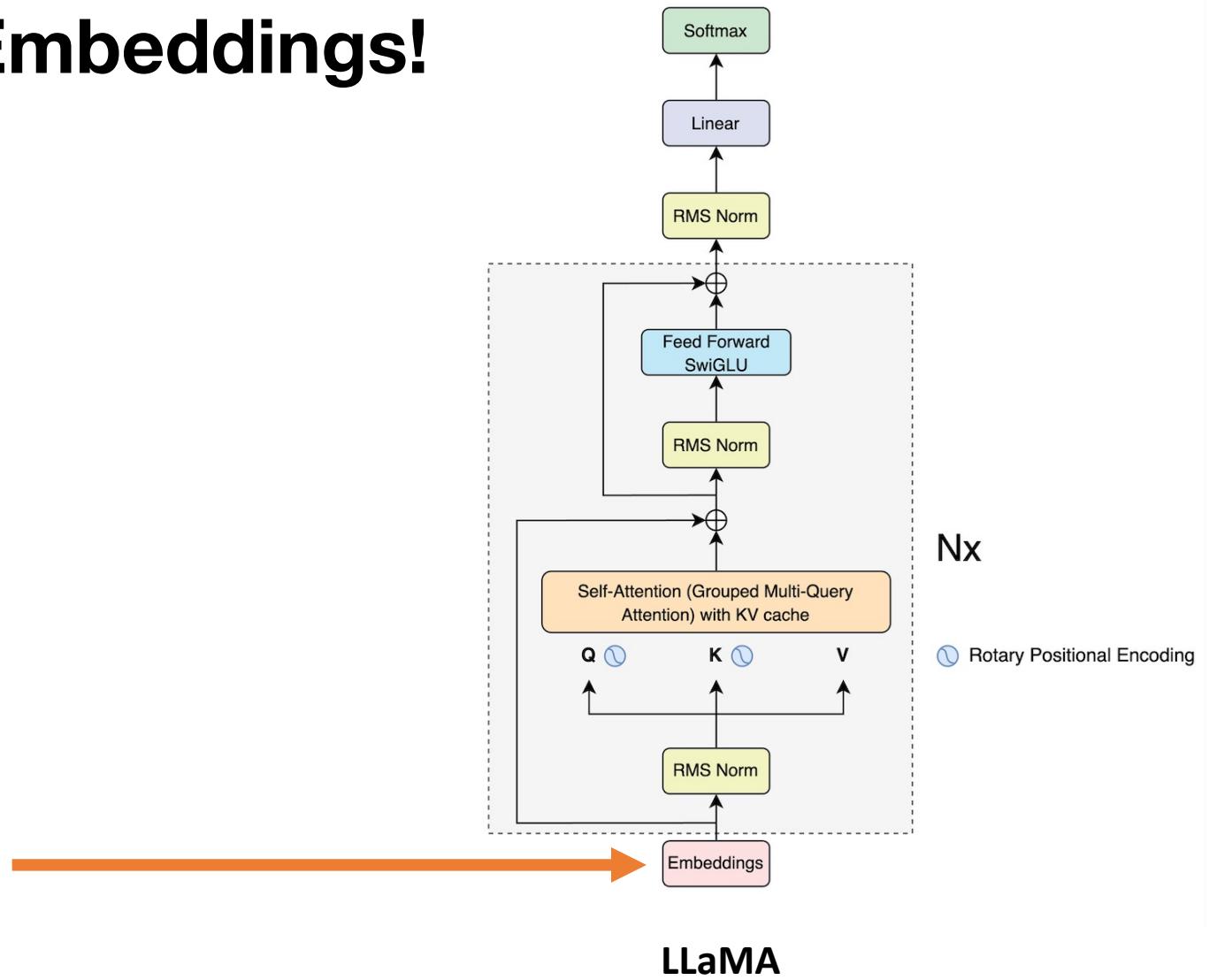
Table 2: **Model sizes, architectures, and optimization hyper-parameters.**

# Models (LLaMA 2)

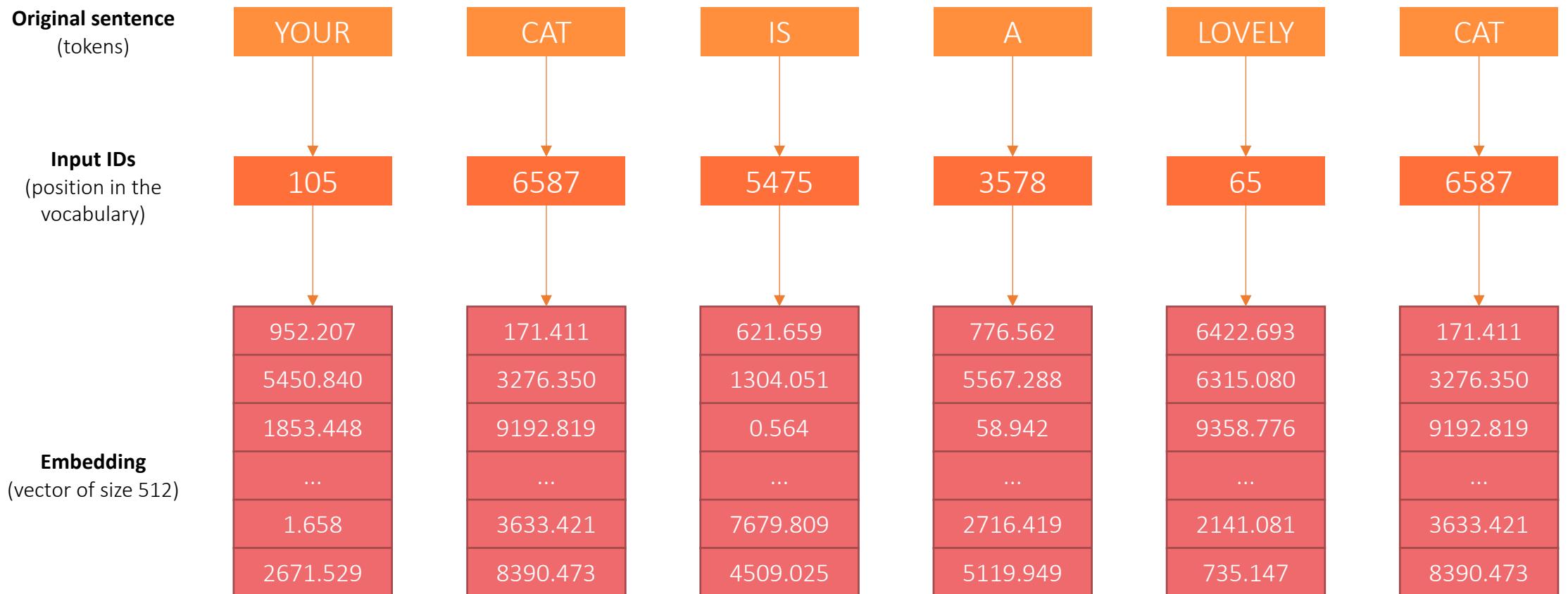
	Training Data	Params	Context Length	GQA	Tokens	LR
LLAMA 1	<i>See Touvron et al. (2023)</i>	7B	2k	✗	1.0T	$3.0 \times 10^{-4}$
		13B	2k	✗	1.0T	$3.0 \times 10^{-4}$
		33B	2k	✗	1.4T	$1.5 \times 10^{-4}$
		65B	2k	✗	1.4T	$1.5 \times 10^{-4}$
LLAMA 2	<i>A new mix of publicly available online data</i>	7B	4k	✗	2.0T	$3.0 \times 10^{-4}$
		13B	4k	✗	2.0T	$3.0 \times 10^{-4}$
		34B	4k	✓	2.0T	$1.5 \times 10^{-4}$
		70B	4k	✓	2.0T	$1.5 \times 10^{-4}$

**Table 1: LLAMA 2 family of models.** Token counts refer to pretraining data only. All models are trained with a global batch-size of 4M tokens. Bigger models — 34B and 70B — use Grouped-Query Attention (GQA) for improved inference scalability.

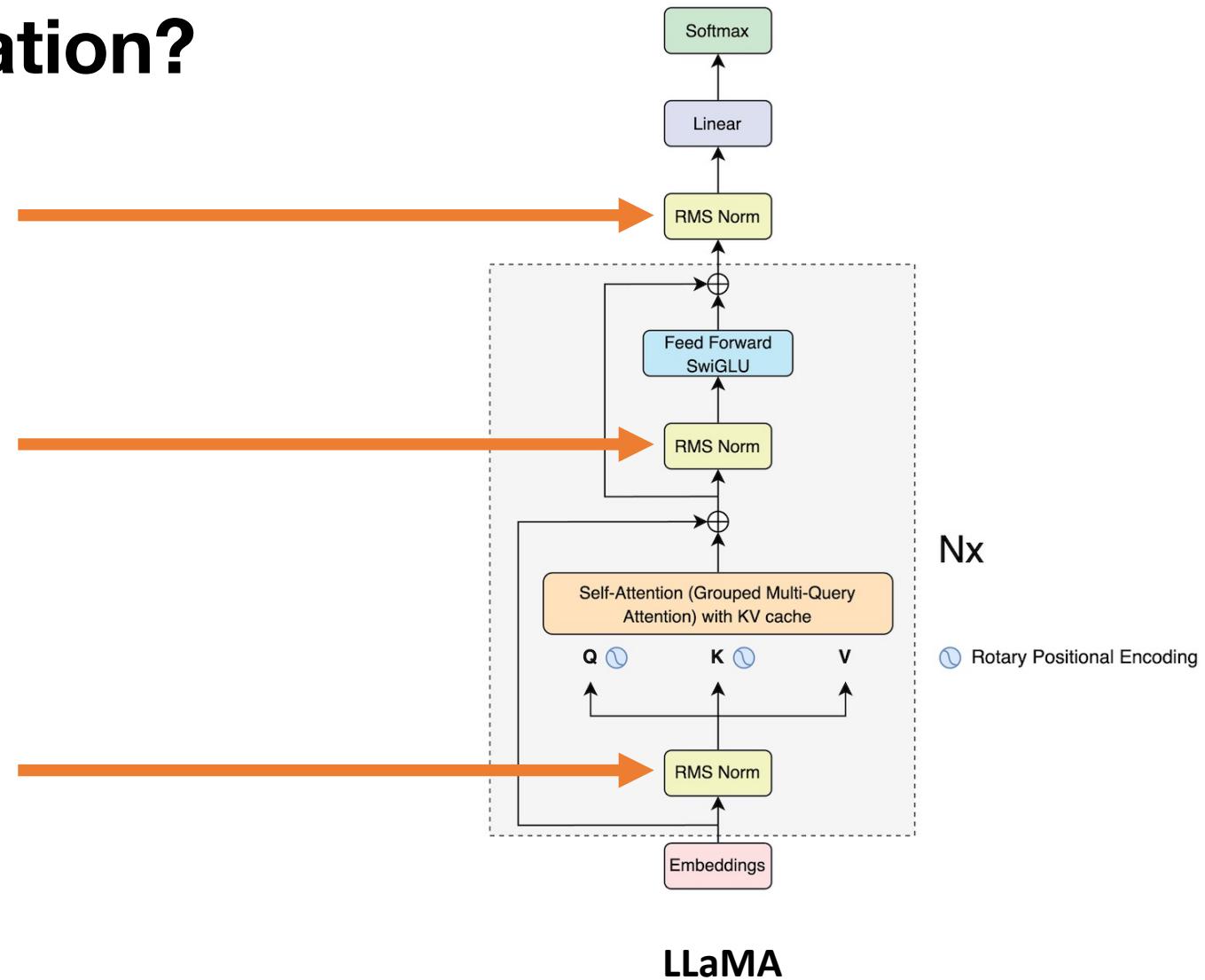
# Let's review the Embeddings!



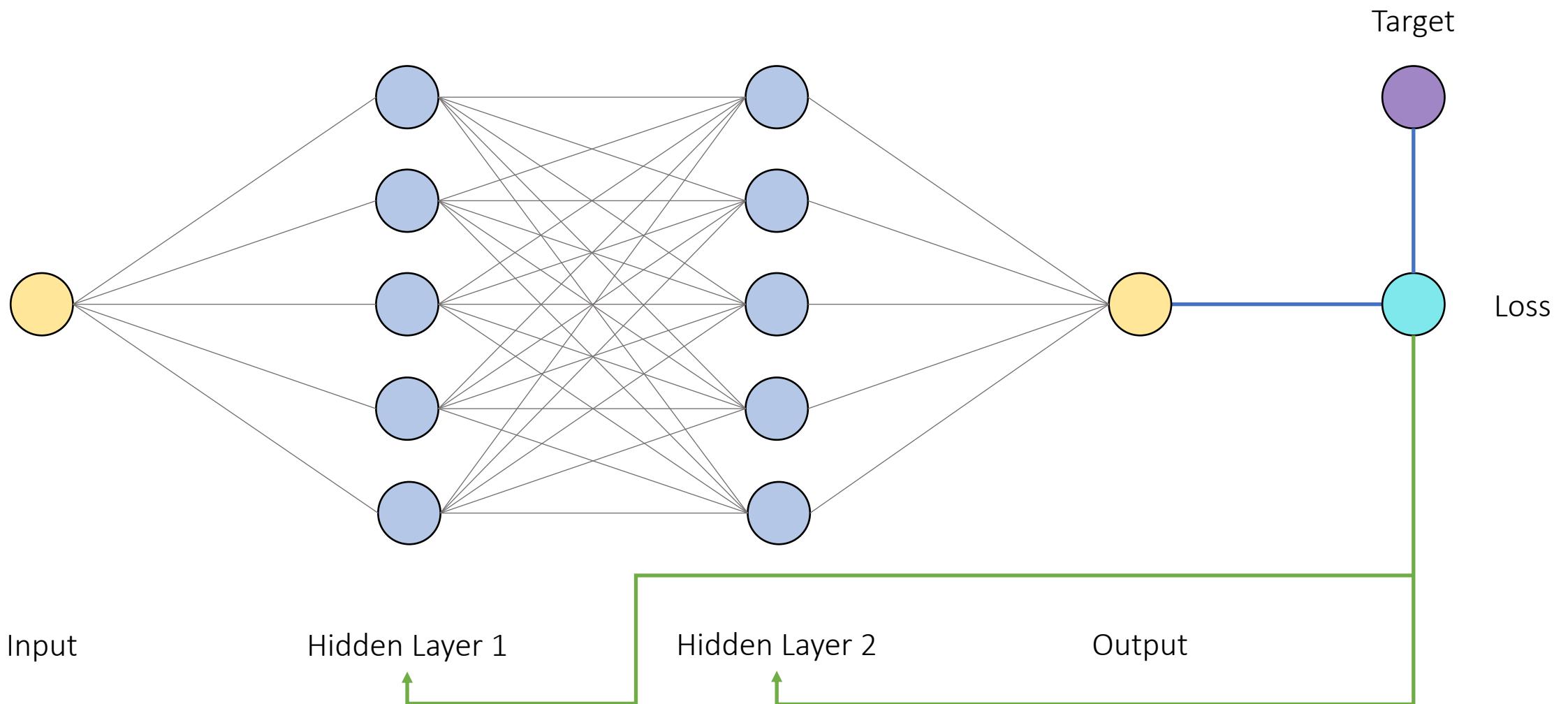
# What is an input embedding?



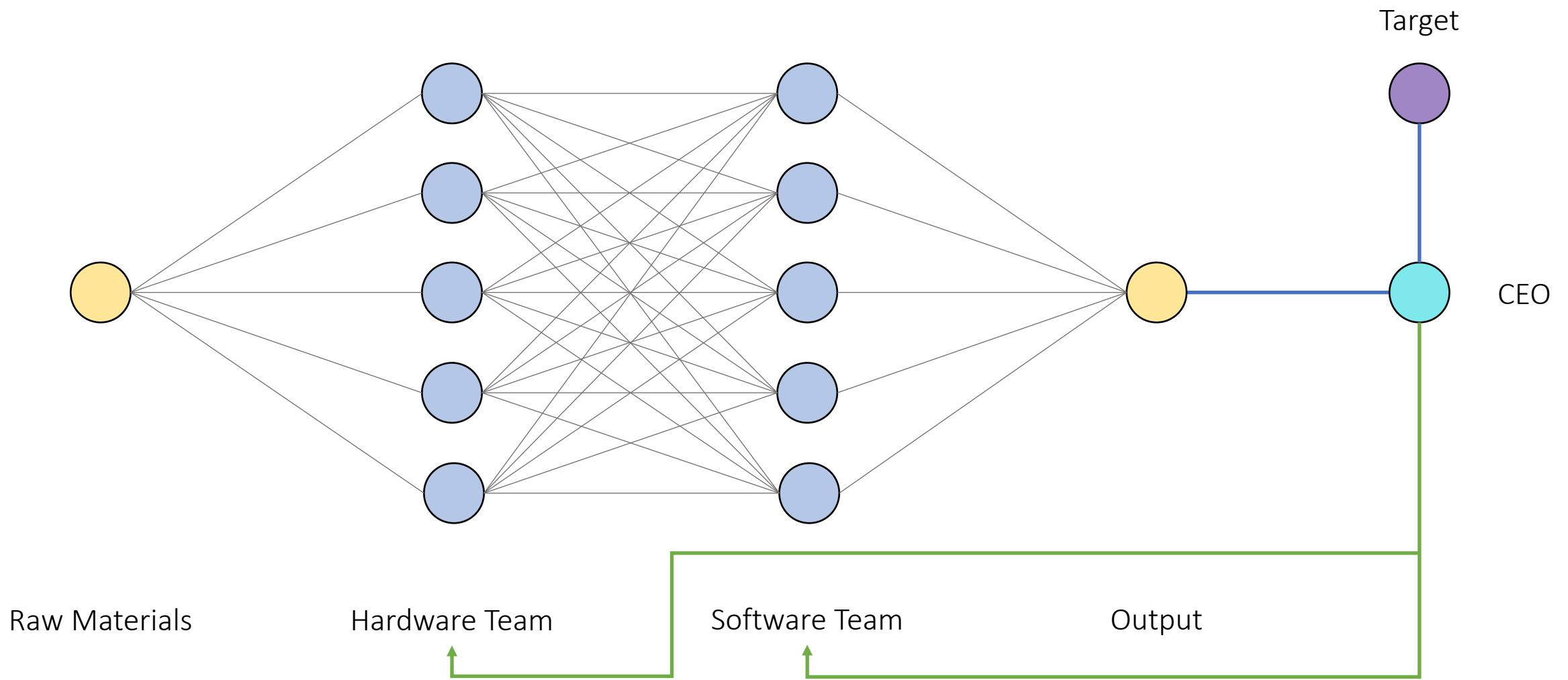
# What is normalization?



# Let's review a simple neural network!



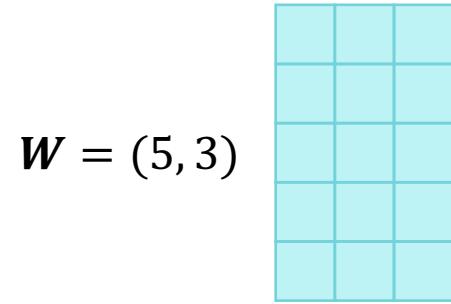
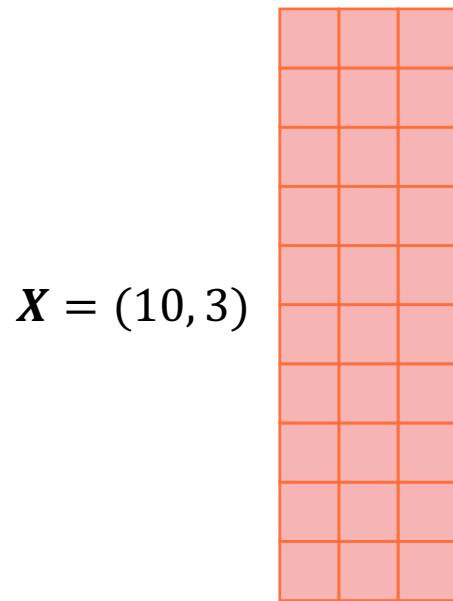
# A simple parallel: The bad CEO in a phone factory



# Let's review the math in neural networks!

Suppose we have a linear layer, defined as `nn.Linear(in_features=3, out_features=5, bias=True)`. This linear layer will create two matrices, called  $W$  (weights) and  $b$  (bias). If we have input  $X$  of shape (10, 3), the output  $O$  will have a shape of (10, 5). But how does this happen mathematically?

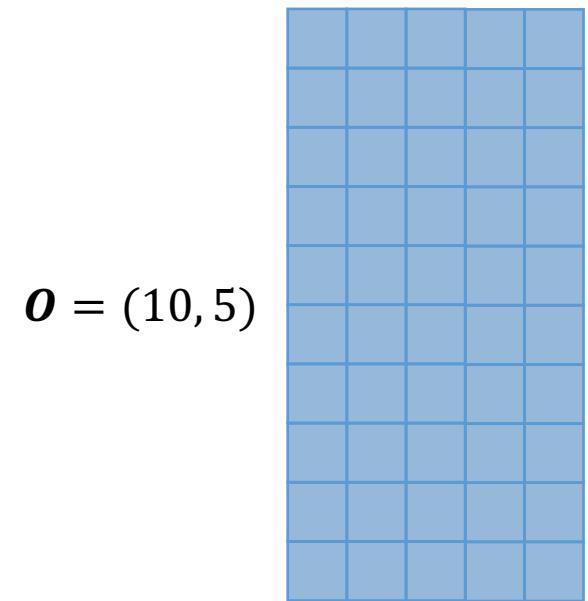
$$O = XW^T + b$$



\*We usually apply a non-linearity  
to the output matrix  $O$

$$b = (1, 5)$$

A horizontal vector of five purple squares, representing the bias vector  $b$  of shape (1, 5).



Each neuron has 3 weights, one for each of the input feature. Each neuron has 1 bias that is added.

# Let's review the math in neural networks!

$$z_1 = (r_1 + b_1) = \left( \sum_{i=1}^3 a_i w_i + b_1 \right)$$

The output of neuron 1 only depends on the first-row features of  $X$ . Usually, we apply a non-linearity like the ReLU function to the output,  $z_1$ .  $z_1$  is referred to as the activation of neuron 1 with respect to the first-row features of  $X$ .

$$\mathbf{b} = (1, 5) \quad \begin{array}{|c|c|c|c|c|} \hline b_1 & & & & \\ \hline \end{array}$$

The bias vector will be broadcasted to every row in the  $XW^T$  table.

$$O = XW^T + b$$

$$W^T = (3, 5)$$

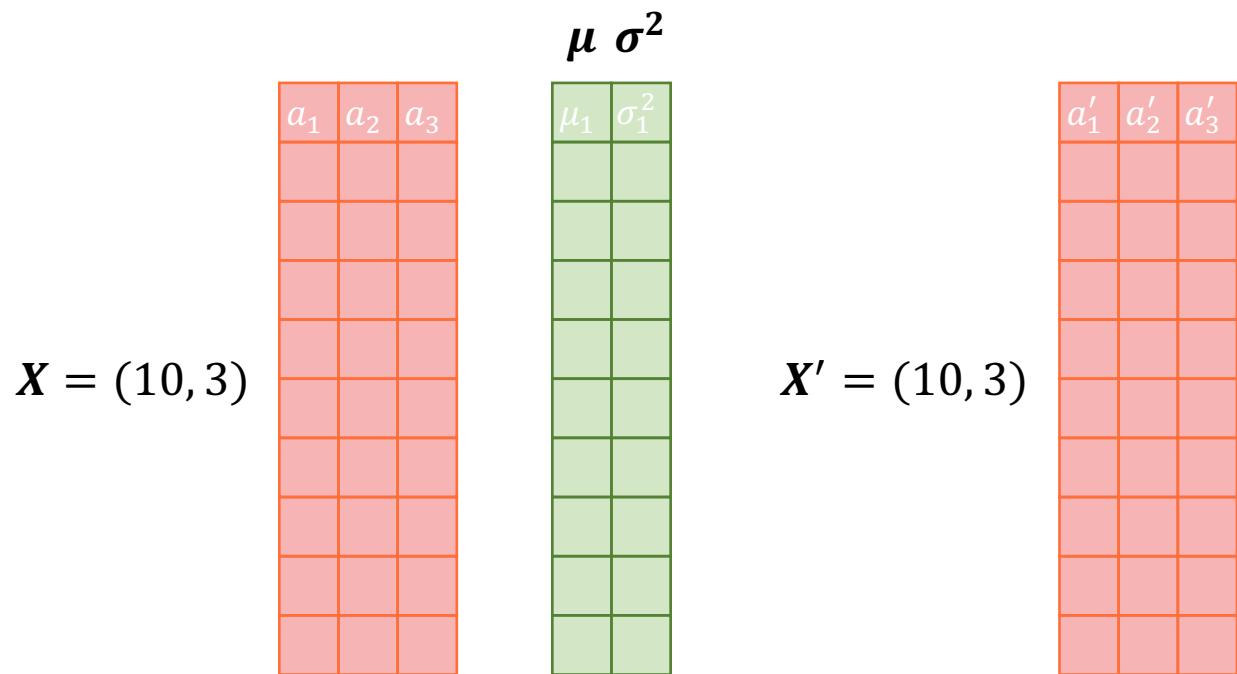
$$XW^T = (10, 5)$$

A 10x5 grid representing a 10x5 matrix  $O = (10, 5)$ . The top-left cell contains the label  $z_1$ .

# Let's review the math in neural networks!

- The output of a neuron depends on the features of the input data (and the neuron's parameters).
- We can think of the input to a neuron as the output of the previous layer.
- If the previous layer (after its weights are updated because of gradient descent) changes drastically, its output (the next layer), will have its input changed drastically as well. So, it will be forced to re-adjust its weights drastically in turn at the next step of gradient descent.
- The phenomenon by which the distributions of internal nodes (neurons) of a neural network change is referred to as **Internal Covariate Shift**. And we want to avoid it because it makes training neural networks slower, as the neurons are forced to drastically re-adjust their weights in one direction or another because of drastic changes in the outputs of the previous layer.

# A Solution: Layer Normalization



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

- Each item is updated with its normalized value, which will turn it into a normal distribution with 0 mean and variance of 1.
- These two parameters **gamma** and **beta** are learnable parameters that allow the model to “amplify” the scale of each feature or apply a translation to the feature according to the needs of the loss function.

With Batch Normalization, we normalize by **columns (features)**.

With Layer Normalization, we normalize by **rows (data items)**.

# Root Mean Square Normalization

---

## Root Mean Square Layer Normalization

---

Biao Zhang<sup>1</sup> Rico Sennrich<sup>2,1</sup>

<sup>1</sup>School of Informatics, University of Edinburgh

<sup>2</sup>Institute of Computational Linguistics, University of Zurich

B.Zhang@ed.ac.uk, sennrich@cl.uzh.ch

## 4 RMSNorm

A well-known explanation of the success of LayerNorm is its re-centering and re-scaling invariance property. The former enables the model to be insensitive to shift noises on both inputs and weights, and the latter keeps the output representations intact when both inputs and weights are randomly scaled. In this paper, we hypothesize that the re-scaling invariance is the reason for success of LayerNorm, rather than re-centering invariance.

We propose RMSNorm which only focuses on re-scaling invariance and regularizes the summed inputs simply according to the root mean square (RMS) statistic:

$$\bar{a}_i = \frac{a_i}{\text{RMS}(\mathbf{a})} g_i, \quad \text{where } \text{RMS}(\mathbf{a}) = \sqrt{\frac{1}{n} \sum_{i=1}^n a_i^2}. \quad (4)$$

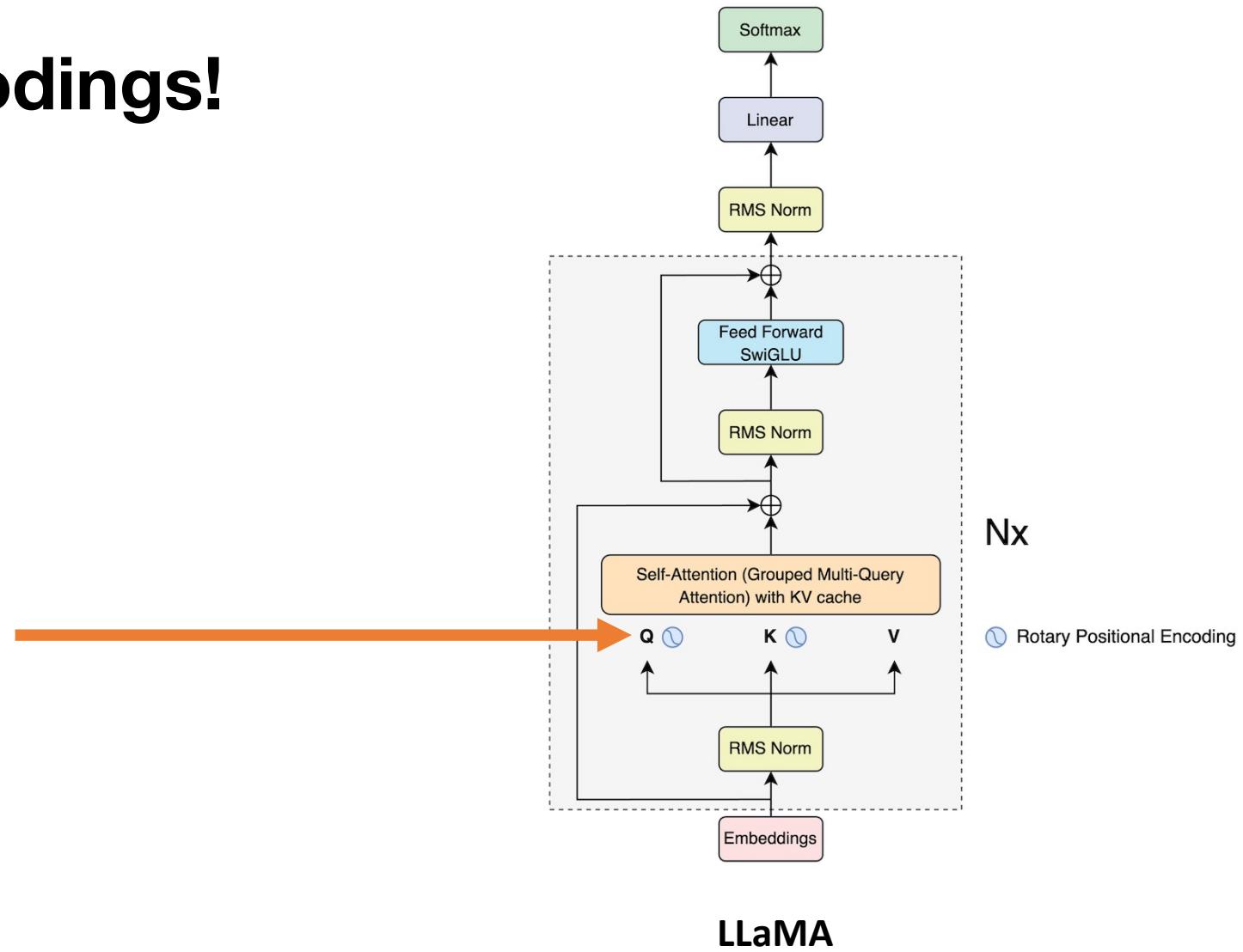
Intuitively, RMSNorm simplifies LayerNorm by totally removing the mean statistic in Eq. (3) at the cost of sacrificing the invariance that mean normalization affords. When the mean of summed inputs is zero, RMSNorm is exactly equal to LayerNorm. Although RMSNorm does not re-center

Just like Layer Normalization, we also have a learnable parameter **gamma** ( $\mathbf{g}$  in the formula on the left) that is multiplied by the normalized values.

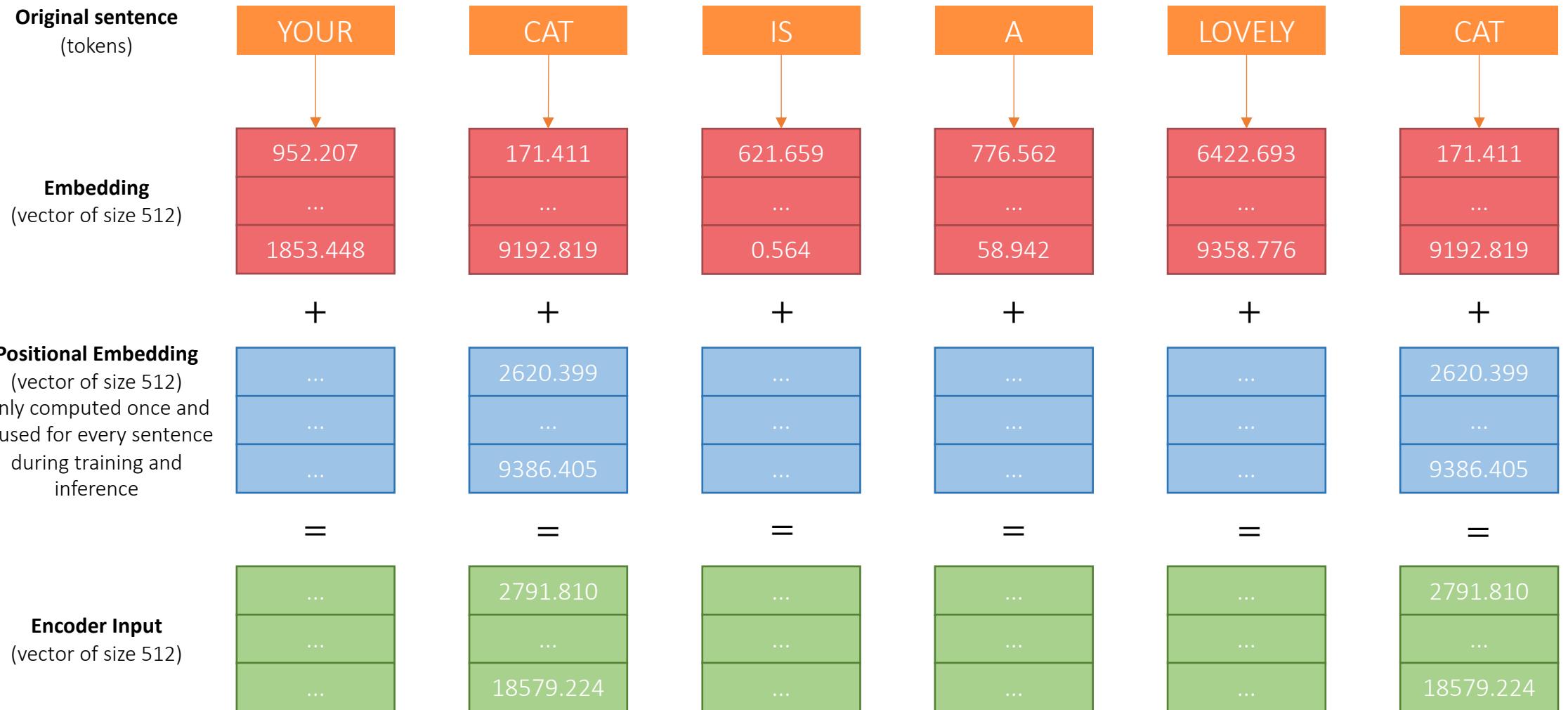
# Why RMSNorm?

- Requires less computation compared to Layer Normalization.
- It works well in practice.

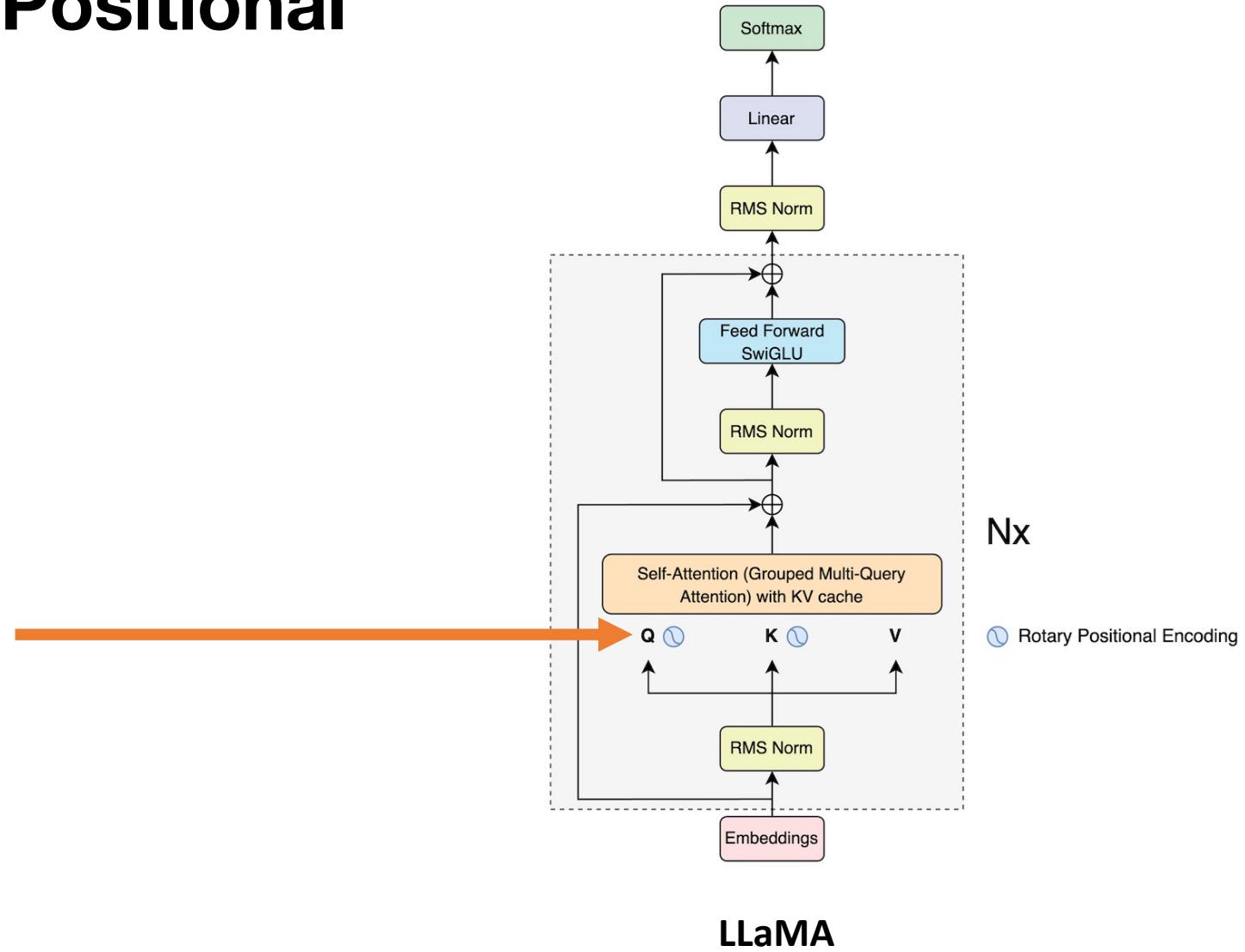
# Let's review Positional Encodings!



# What is an input embedding?



# What is Rotary Positional Encoding?



# What's the difference between the absolute positional encodings and the relative ones?

- Absolute positional encodings are fixed vectors that are added to the embedding of a token to represent its absolute position in the sentence. So, it deals with **one token at a time**. You can think of it as the latitude and longitude on a map: each point on earth will have a unique pair.
- Relative positional encodings, on the other hand, deals with two tokens at a time and it is involved when we calculate the attention: since the attention mechanism captures the “intensity” of how much the two words are related to each other. Relative positional encodings tell the attention mechanism the **distance** between the two words involved in it. So, given **two tokens**, we create a vector that represents their distance.
- Relative positional encodings were introduced in the following paper:

## Self-Attention with Relative Position Representations

**Peter Shaw**

Google

[petershaw@google.com](mailto:petershaw@google.com)

**Jakob Uszkoreit**

Google Brain

[usz@google.com](mailto:usz@google.com)

**Ashish Vaswani**

Google Brain

[avaswani@google.com](mailto:avaswani@google.com)

**Absolute** Positional Encodings  
From “Attention is All You Need”

$$e_{ij} = \frac{(x_i W^Q)(x_j W^K)^T}{\sqrt{d_z}}$$



**Relative** Positional Encodings  
From “Self-Attention with Relative Position Representations”

$$e_{ij} = \frac{x_i W^Q (x_j W^K + a_{ij}^K)^T}{\sqrt{d_z}}$$



\*NSFW: Menage a trois

# Rotary Positional Embeddings

---

## RoFORMER: ENHANCED TRANSFORMER WITH ROTARY POSITION EMBEDDING

---

**Jianlin Su**

Zhuiyi Technology Co., Ltd.  
Shenzhen  
[bojonesu@wezhuiyi.com](mailto:bojonesu@wezhuiyi.com)

**Yu Lu**

Zhuiyi Technology Co., Ltd.  
Shenzhen  
[julianlu@wezhuiyi.com](mailto:julianlu@wezhuiyi.com)

**Shengfeng Pan**

Zhuiyi Technology Co., Ltd.  
Shenzhen  
[nickpan@wezhuiyi.com](mailto:nickpan@wezhuiyi.com)

**Ahmed Murtadha**

Zhuiyi Technology Co., Ltd.  
Shenzhen  
[mengjiayi@wezhuiyi.com](mailto:mengjiayi@wezhuiyi.com)

**Bo Wen**

Zhuiyi Technology Co., Ltd.  
Shenzhen  
[brucewen@wezhuiyi.com](mailto:brucewen@wezhuiyi.com)

**Yunfeng Liu**

Zhuiyi Technology Co., Ltd.  
Shenzhen  
[glenliu@wezhuiyi.com](mailto:glenliu@wezhuiyi.com)

# Rotary Position Embeddings: The inner product

- The **dot product** used in the attention mechanism is a type of *inner* product, which can be thought of as a generalization of the dot product.
- Can we find an inner product over the two vectors  $\mathbf{q}$  (query) and  $\mathbf{k}$  (key) used in the attention mechanism that only depends on the two vectors and the relative distance of the token they represent?

Under the case of  $d = 2$ , we consider two-word embedding vectors  $\mathbf{x}_q, \mathbf{x}_k$  corresponds to query and key and their position  $m$  and  $n$ , respectively. According to eq. (1), their position-encoded counterparts are:

$$\begin{aligned}\mathbf{q}_m &= f_q(\mathbf{x}_q, m), \\ \mathbf{k}_n &= f_k(\mathbf{x}_k, n),\end{aligned}\tag{20}$$

where the subscripts of  $\mathbf{q}_m$  and  $\mathbf{k}_n$  indicate the encoded positions information. Assume that there exists a function  $g$  that defines the inner product between vectors produced by  $f_{\{q,k\}}$ :

$$\mathbf{q}_m^\top \mathbf{k}_n = \langle f_q(\mathbf{x}_m, m), f_k(\mathbf{x}_n, n) \rangle = g(\mathbf{x}_m, \mathbf{x}_n, n - m),\tag{21}$$

# Rotary Position Embeddings: The inner product

- We can define a function  $g$  like the following that only depends on the two embedding vectors,  $\mathbf{q}$  and  $\mathbf{k}$  and their relative distance.

$$f_q(\mathbf{x}_m, m) = (\mathbf{W}_q \mathbf{x}_m) e^{im\theta}$$

$$f_k(\mathbf{x}_n, n) = (\mathbf{W}_k \mathbf{x}_n) e^{in\theta}$$

$$g(\mathbf{x}_m, \mathbf{x}_n, m - n) = \text{Re}[(\mathbf{W}_q \mathbf{x}_m)(\mathbf{W}_k \mathbf{x}_n)^* e^{i(m-n)\theta}]$$

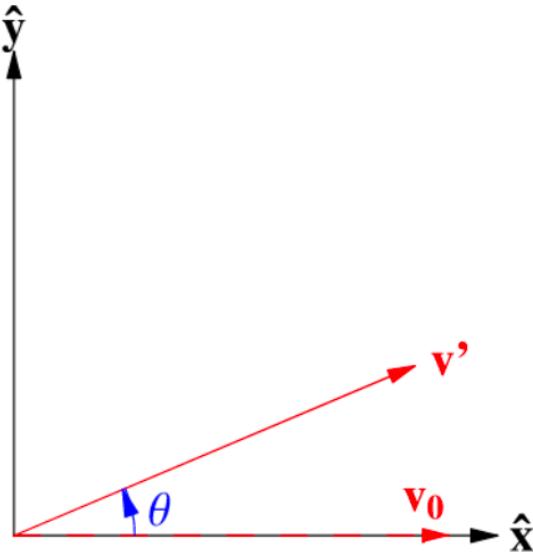
\* = **Conjugate** of the complex number

- Using **Euler's formula**, we can write it into its matrix form.

$$f_{\{q,k\}}(\mathbf{x}_m, m) = \begin{pmatrix} \cos m\theta & -\sin m\theta \\ \sin m\theta & \cos m\theta \end{pmatrix} \begin{pmatrix} W_{\{q,k\}}^{(11)} & W_{\{q,k\}}^{(12)} \\ W_{\{q,k\}}^{(21)} & W_{\{q,k\}}^{(22)} \end{pmatrix} \begin{pmatrix} x_m^{(1)} \\ x_m^{(2)} \end{pmatrix}$$

Rotation matrix in a 2D space, hence the name **rotary** position embeddings

# Rotary Position Embeddings: The rotation matrix



In  $\mathbb{R}^2$ , consider the matrix that rotates a given vector  $v_0$  by a counterclockwise angle  $\theta$  in a fixed coordinate system. Then

$$R_\theta = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}, \quad (1)$$

so

$$v' = R_\theta v_0. \quad (2)$$

# Rotary Position Embeddings: The general form

Since the matrix is **sparse**, it is not convenient to use it to compute the positional embeddings.

$$f_{\{q,k\}}(\mathbf{x}_m, m) = \mathbf{R}_{\Theta,m}^d \mathbf{W}_{\{q,k\}} \mathbf{x}_m \quad (14)$$

where

$$\mathbf{R}_{\Theta,m}^d = \begin{pmatrix} \cos m\theta_1 & -\sin m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ \sin m\theta_1 & \cos m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cos m\theta_2 & -\sin m\theta_2 & \cdots & 0 & 0 \\ 0 & 0 & \sin m\theta_2 & \cos m\theta_2 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & \cos m\theta_{d/2} & -\sin m\theta_{d/2} \\ 0 & 0 & 0 & 0 & \cdots & \sin m\theta_{d/2} & \cos m\theta_{d/2} \end{pmatrix} \quad (15)$$

is the rotary matrix with pre-defined parameters  $\Theta = \{\theta_i = 10000^{-2(i-1)/d}, i \in [1, 2, \dots, d/2]\}$ . A graphic illustration of RoPE is shown in Figure (1). Applying our RoPE to self-attention in Equation (2), we obtain:

$$\mathbf{q}_m^\top \mathbf{k}_n = (\mathbf{R}_{\Theta,m}^d \mathbf{W}_q \mathbf{x}_m)^\top (\mathbf{R}_{\Theta,n}^d \mathbf{W}_k \mathbf{x}_n) = \mathbf{x}_m^\top \mathbf{W}_q \mathbf{R}_{\Theta,n-m}^d \mathbf{W}_k \mathbf{x}_n \quad (16)$$

where  $\mathbf{R}_{\Theta,n-m}^d = (\mathbf{R}_{\Theta,m}^d)^\top \mathbf{R}_{\Theta,n}^d$ . Note that  $\mathbf{R}_{\Theta}^d$  is an orthogonal matrix, which ensures stability during the process of encoding position information. In addition, due to the sparsity of  $R_{\Theta}^d$ , applying matrix multiplication directly as in Equation (16) is not computationally efficient; we provide another realization in theoretical explanation.

# Rotary Position Embeddings: The computationally-efficient form

Given a token with embedding vector  $\mathbf{x}$ , and the position  $\mathbf{m}$  of the token inside the sentence, this is how we compute the position embeddings for the token.

## 3.4.2 Computational efficient realization of rotary matrix multiplication

Taking the advantage of the sparsity of  $R_{\Theta,m}^d$  in Equation (15), a more computational efficient realization of a multiplication of  $R_\Theta^d$  and  $\mathbf{x} \in \mathbb{R}^d$  is:

$$\mathbf{R}_{\Theta,m}^d \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ \vdots \\ x_{d-1} \\ x_d \end{pmatrix} \otimes \begin{pmatrix} \cos m\theta_1 \\ \cos m\theta_1 \\ \cos m\theta_2 \\ \cos m\theta_2 \\ \vdots \\ \cos m\theta_{d/2} \\ \cos m\theta_{d/2} \end{pmatrix} + \begin{pmatrix} -x_2 \\ x_1 \\ -x_4 \\ x_3 \\ \vdots \\ -x_{d-1} \\ x_d \end{pmatrix} \otimes \begin{pmatrix} \sin m\theta_1 \\ \sin m\theta_1 \\ \sin m\theta_2 \\ \sin m\theta_2 \\ \vdots \\ \sin m\theta_{d/2} \\ \sin m\theta_{d/2} \end{pmatrix} \quad (34)$$

# Rotary Position Embeddings: Long-term decay

The authors calculated an **upper bound** for the inner product by varying the distance between two tokens and proved that it decays with the growth of the relative distance.

This means that the “intensity” of relationship between two tokens encoded with Rotary Positional Embeddings will be numerically smaller as the distance between them grows.

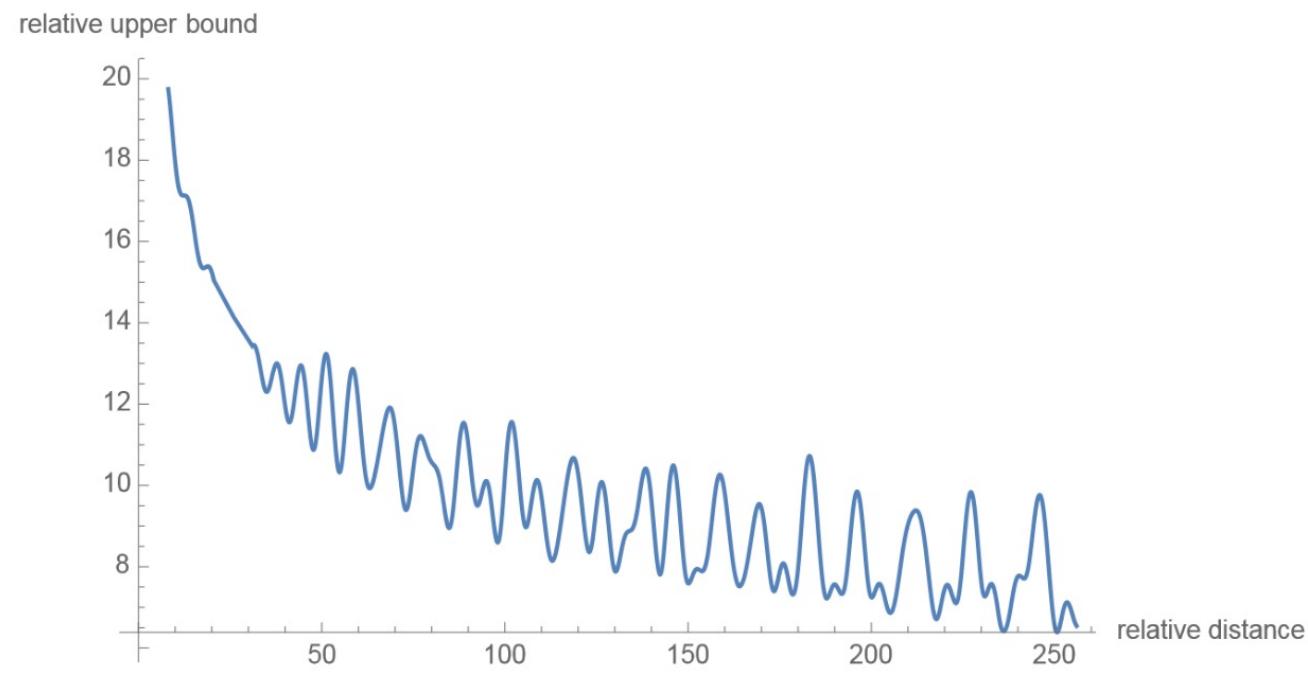
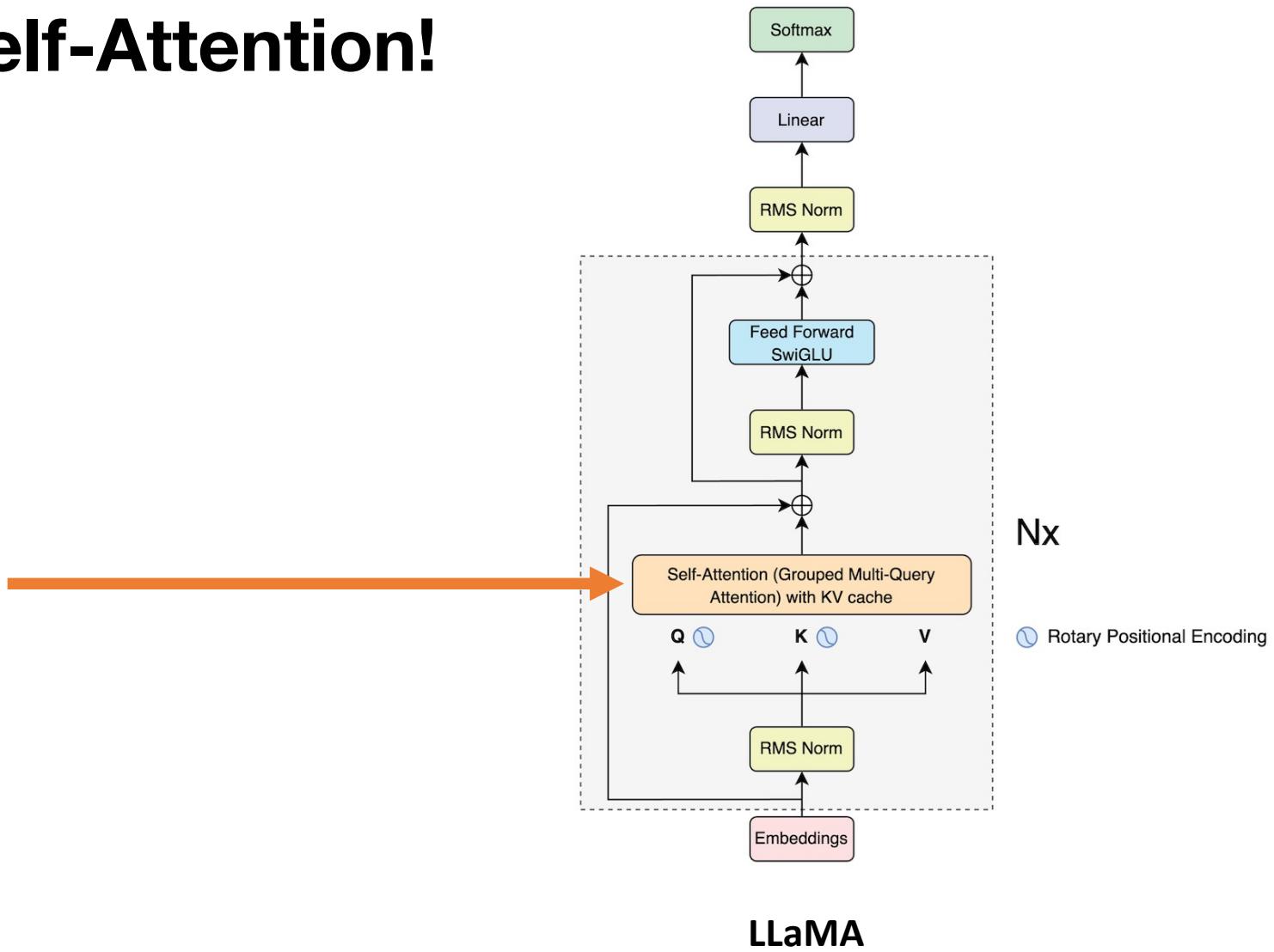


Figure 2: Long-term decay of RoPE.

# Rotary Position Embeddings: Practical considerations

- The rotary position embeddings are **only applied to the query and the keys**, but not the values.
- The rotary position embeddings are applied after the vectors  $\mathbf{q}$  and  $\mathbf{k}$  have been multiplied by the  $\mathbf{W}$  matrix in the attention mechanism, while in the vanilla transformer, they are applied before.

# Let's review Self-Attention!



# What is Self-Attention?

Self-Attention allows the model to relate words to one another.

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

In this simple case, we consider the sequence length  $\text{seq} = 6$  and  $d_{\text{model}} = d_k = 512$ .

The matrices  $Q$ ,  $K$  and  $V$  are just the input sentence.



\* for simplicity, only one head is considered where  $d_{\text{model}} = d_k$

	THE	CAT	IS	ON	A	CHAIR
THE	0.268	0.119	0.134	0.148	0.179	0.152
CAT	0.124	0.278	0.201	0.128	0.154	0.115
IS	0.147	0.132	0.262	0.097	0.218	0.145
ON	0.210	0.128	0.206	0.212	0.119	0.125
A	0.146	0.158	0.152	0.143	0.227	0.174
CHAIR	0.195	0.114	0.203	0.103	0.157	0.229

\* all values are random

(6, 6)

# How to compute Self-Attention?

	THE	CAT	IS	ON	A	CHAIR
THE	0.268	0.119	0.134	0.148	0.179	0.152
CAT	0.124	0.278	0.201	0.128	0.154	0.115
IS	0.147	0.132	0.262	0.097	0.218	0.145
ON	0.210	0.128	0.206	0.212	0.119	0.125
A	0.146	0.158	0.152	0.143	0.227	0.174
CHAIR	0.195	0.114	0.203	0.103	0.157	0.229

(6, 6)

X



=



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

Each row in this matrix captures not only the meaning (given by the embedding) or the position in the sentence (represented by the positional embeddings) but also each word's interaction with other words.

# Next Token Prediction Task

- Imagine we want to train a model to write Dante Alighieri's Divine Comedy's 5<sup>th</sup> Canto from the Inferno.

**Amor, ch'al cor gentil ratto s'apprende,**  
prese costui de la bella persona  
che mi fu tolta; e 'l modo ancor m'offende.

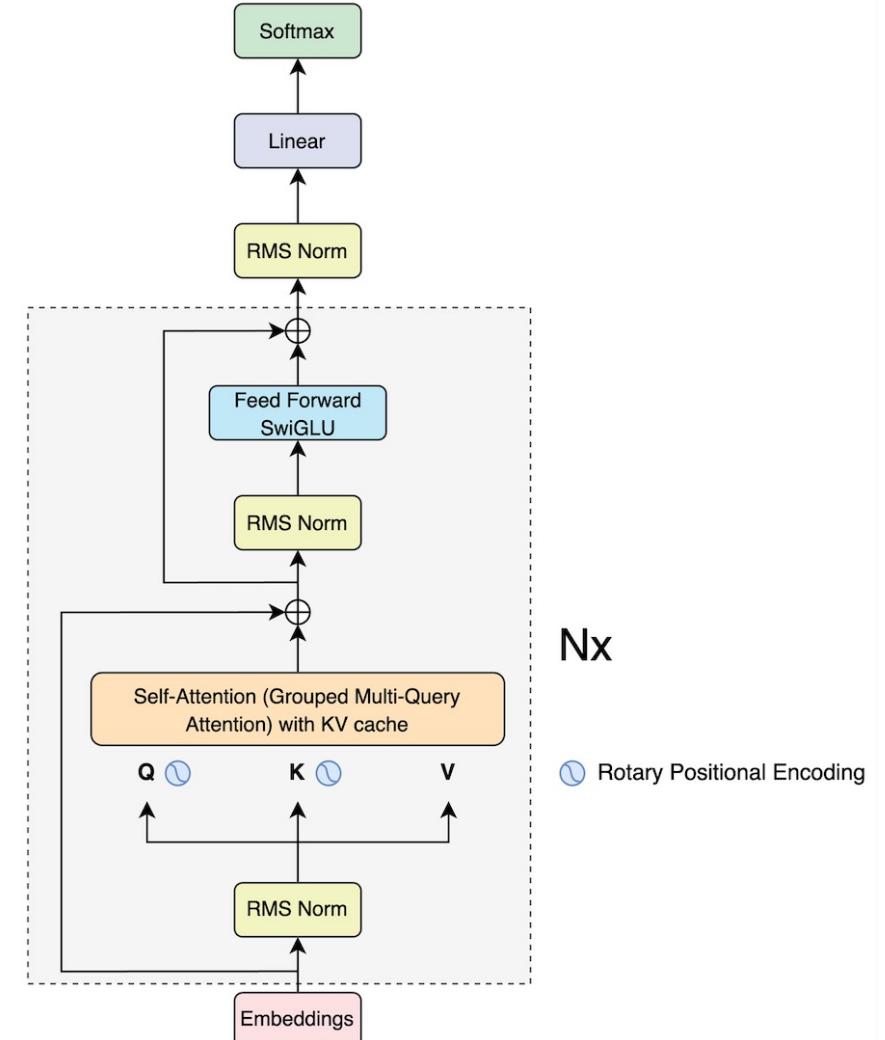
Amor, ch'a nullo amato amar perdona,  
mi prese del costui piacer si forte,  
che, come vedi, ancor non m'abbandona.

Amor condusse noi ad una morte.  
Caina attende chi a vita ci spense.

**Love, that can quickly seize the gentle heart,**  
took hold of him because of the fair body  
taken from me-how that was done still wounds me.

Love, that releases no beloved from loving,  
took hold of me so strongly through his beauty  
that, as you see, it has not left me yet.

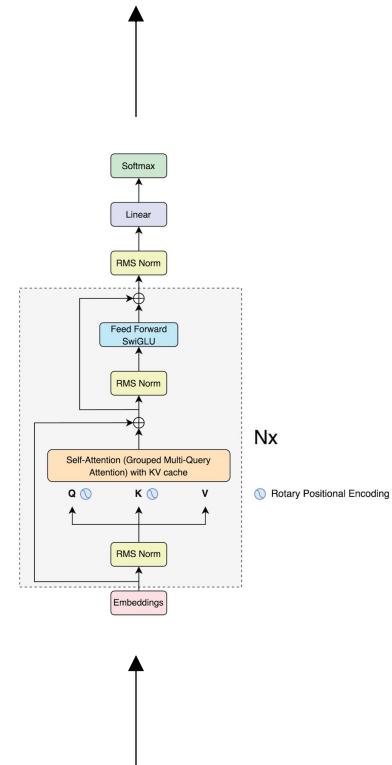
Love led the two of us unto one death.  
Caina waits for him who took our life."



# Next Token Prediction Task

Training

Target Love that can quickly seize the gentle heart [EOS]

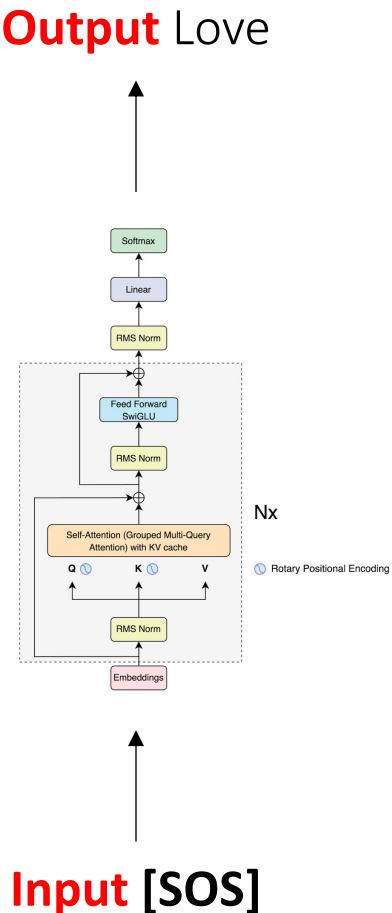


Input [SOS] Love that can quickly seize the gentle heart

# Next Token Prediction Task

Inference

T=1

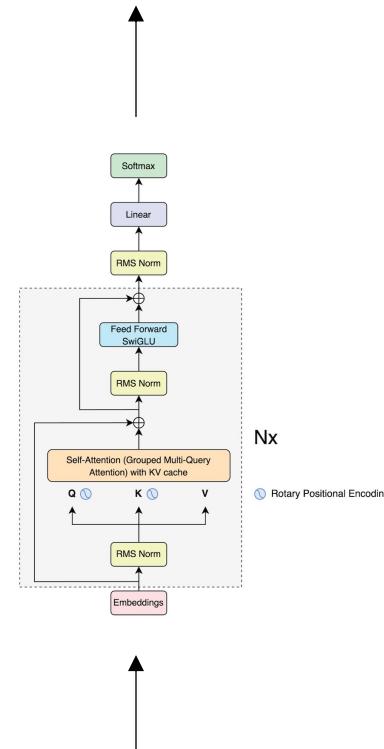


# Next Token Prediction Task

Inference

T=2

**Output** Love that



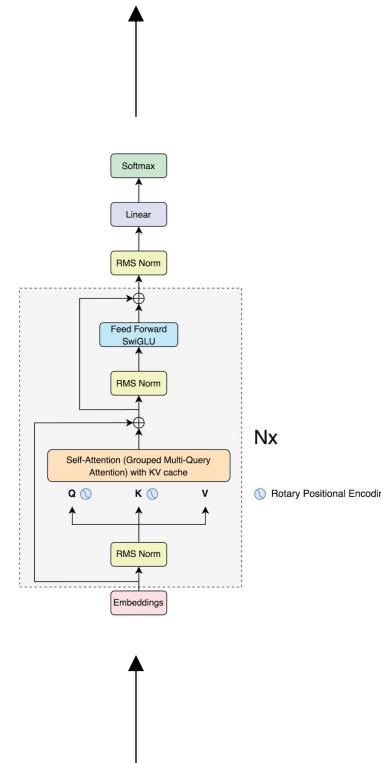
**Input** [SOS] Love

# Next Token Prediction Task

Inference

T=3

**Output** Love that can



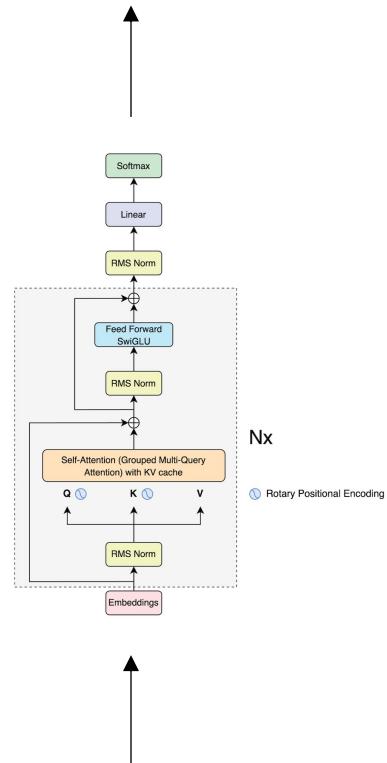
**Input [SOS]** Love that

# Next Token Prediction Task

Inference

T=4

**Output** Love that can quickly



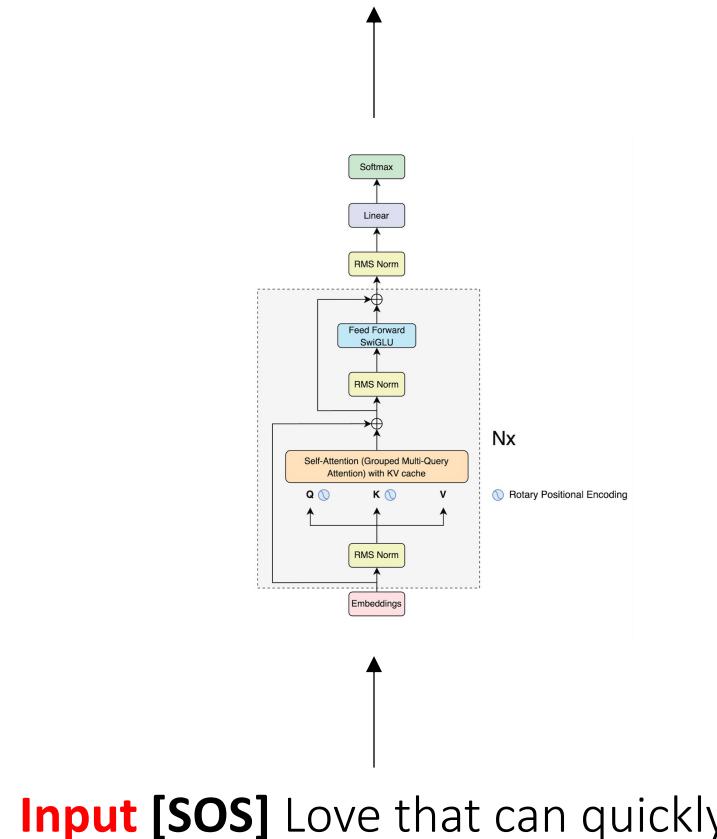
**Input [SOS]** Love that can

# Next Token Prediction Task

Inference

T=5

**Output** Love that can quickly seize



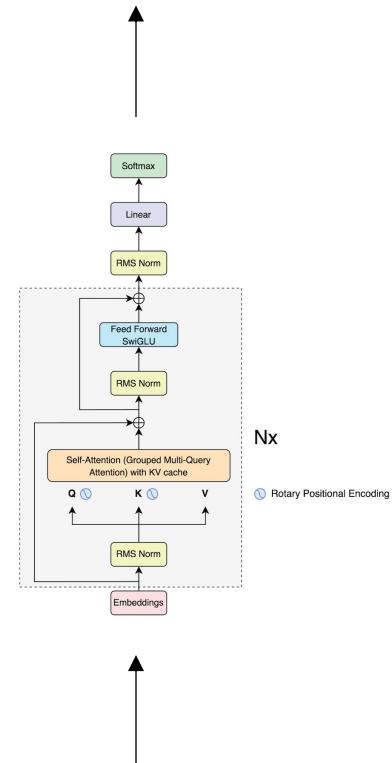
# Next Token Prediction Task

Inference

T=6

**Output** Love that can quickly seize the

**Input [SOS]** Love that can quickly seize

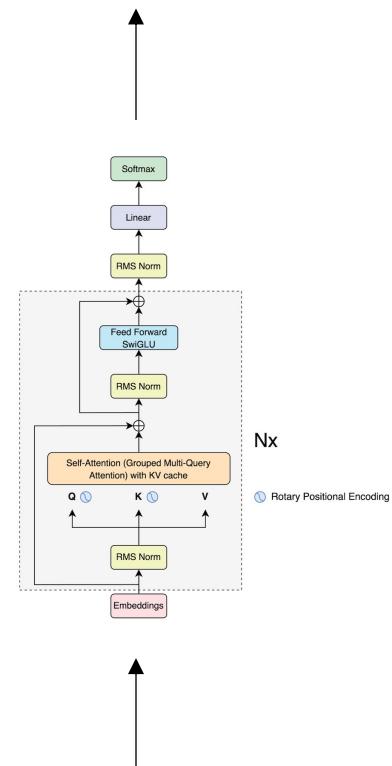


# Next Token Prediction Task

Inference

T=7

**Output** Love that can quickly seize the gentle



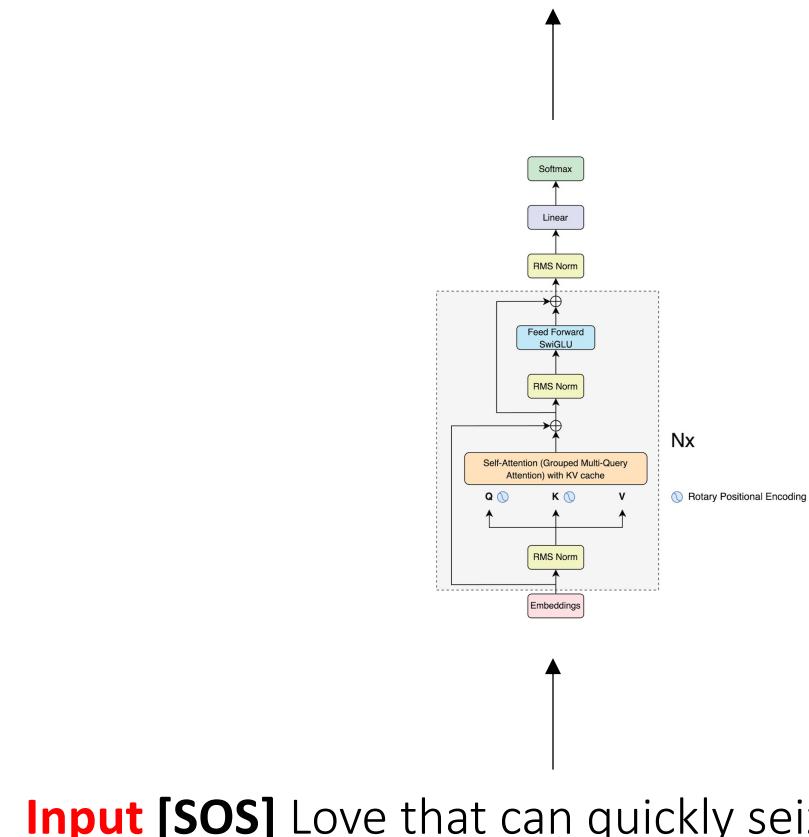
**Input [SOS]** Love that can quickly seize the

# Next Token Prediction Task

Inference

T=8

**Output** Love that can quickly seize the gentle heart

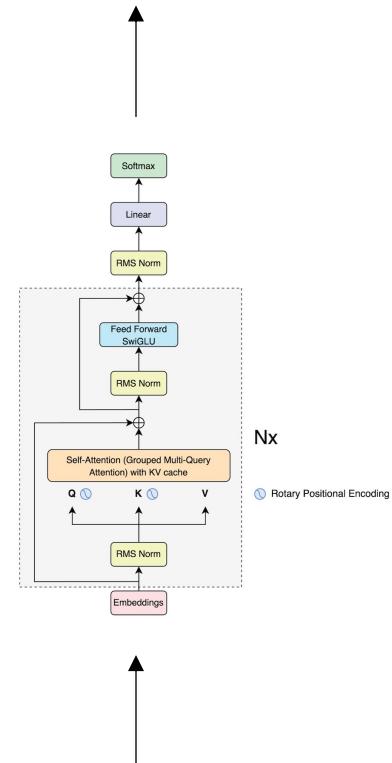


# Next Token Prediction Task

Inference

T=9

**Output** Love that can quickly seize the gentle heart [EOS]

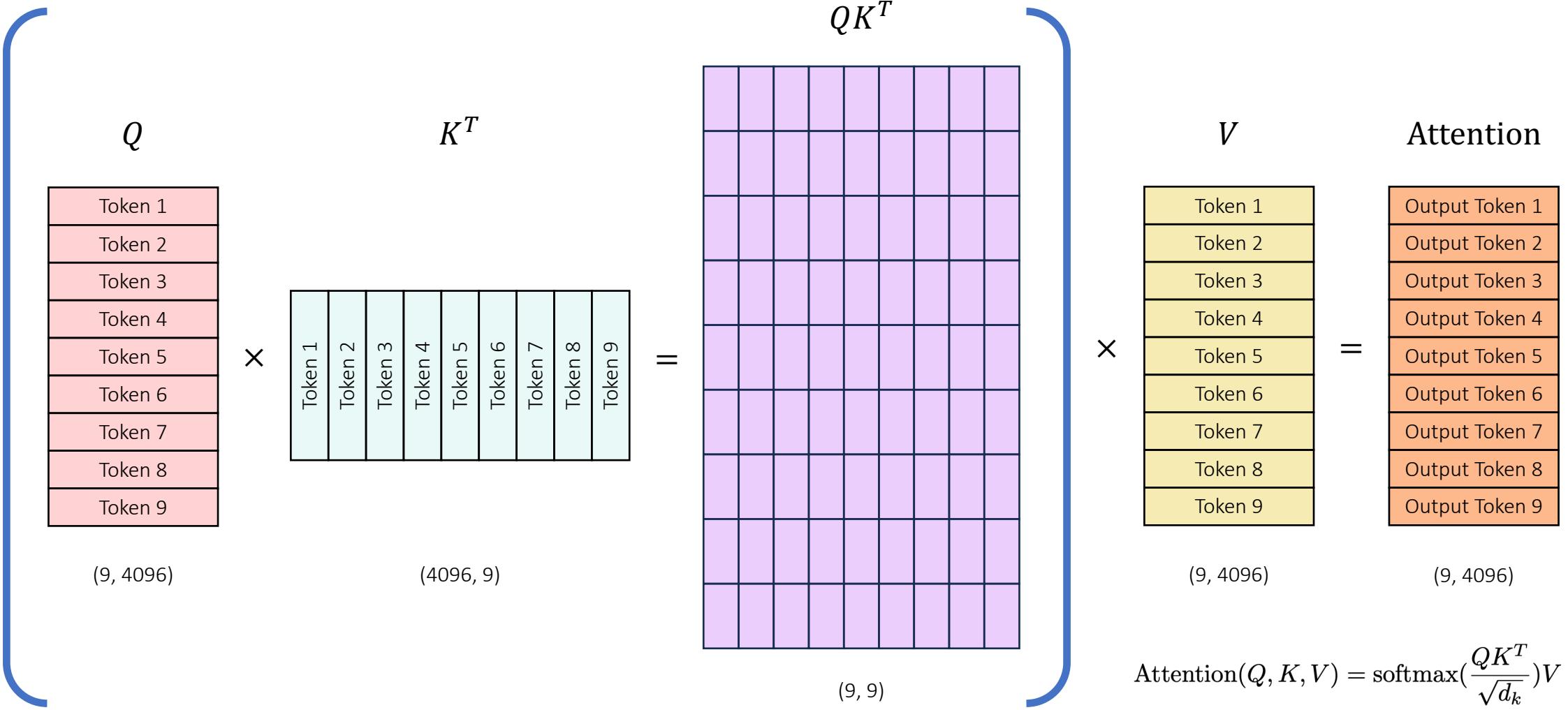


**Input [SOS]** Love that can quickly seize the gentle heart

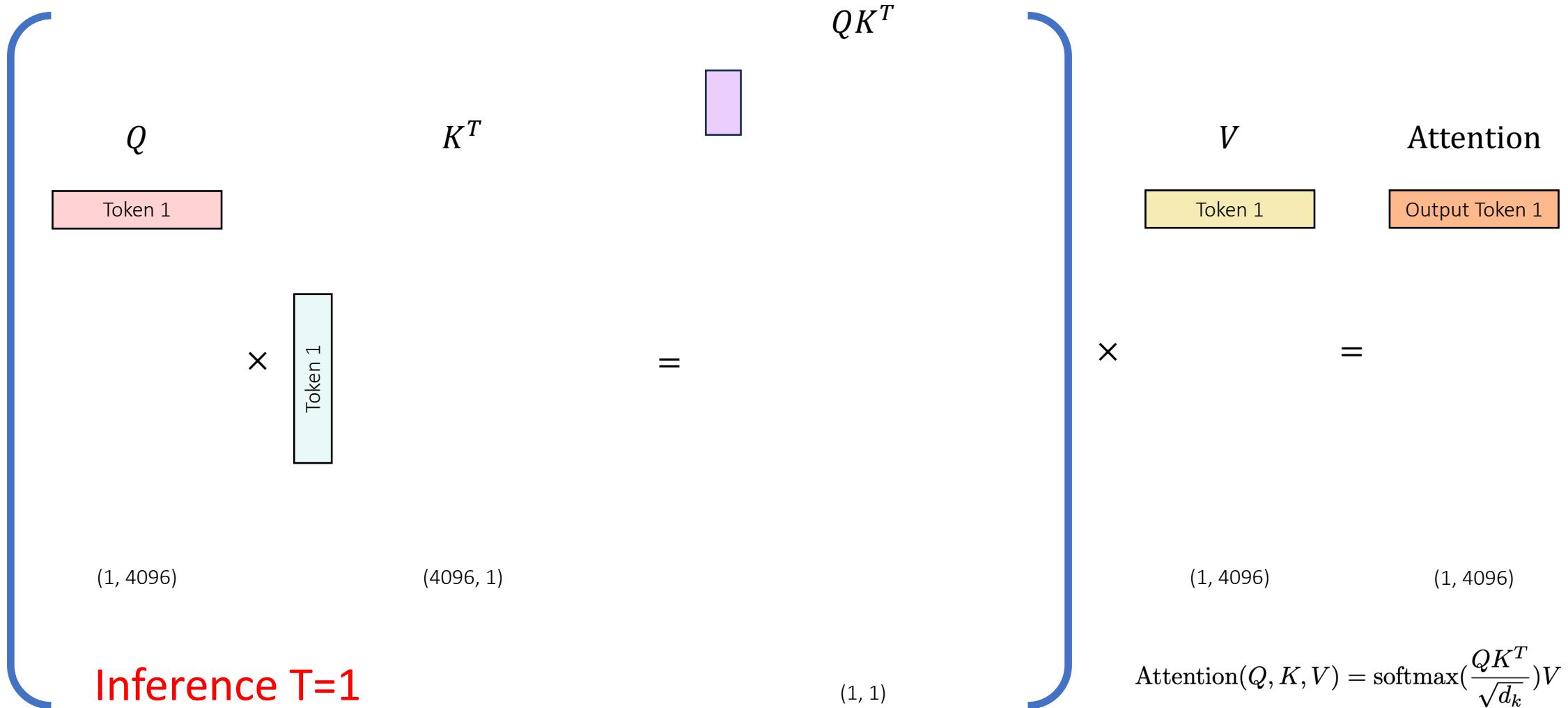
# Next Token Prediction Task: The Motivation behind KV-Cache

- At every step of the inference, we are only interested in the **last token** output by the model, because we already have the previous tokens. However, the model needs access to all previous tokens to decide on which token to output, since they constitute its context (or the “prompt”).
- Is there a way to make the model do less computation on the token it has already seen **during inference**? YES! The solution is **KV-Cache**!

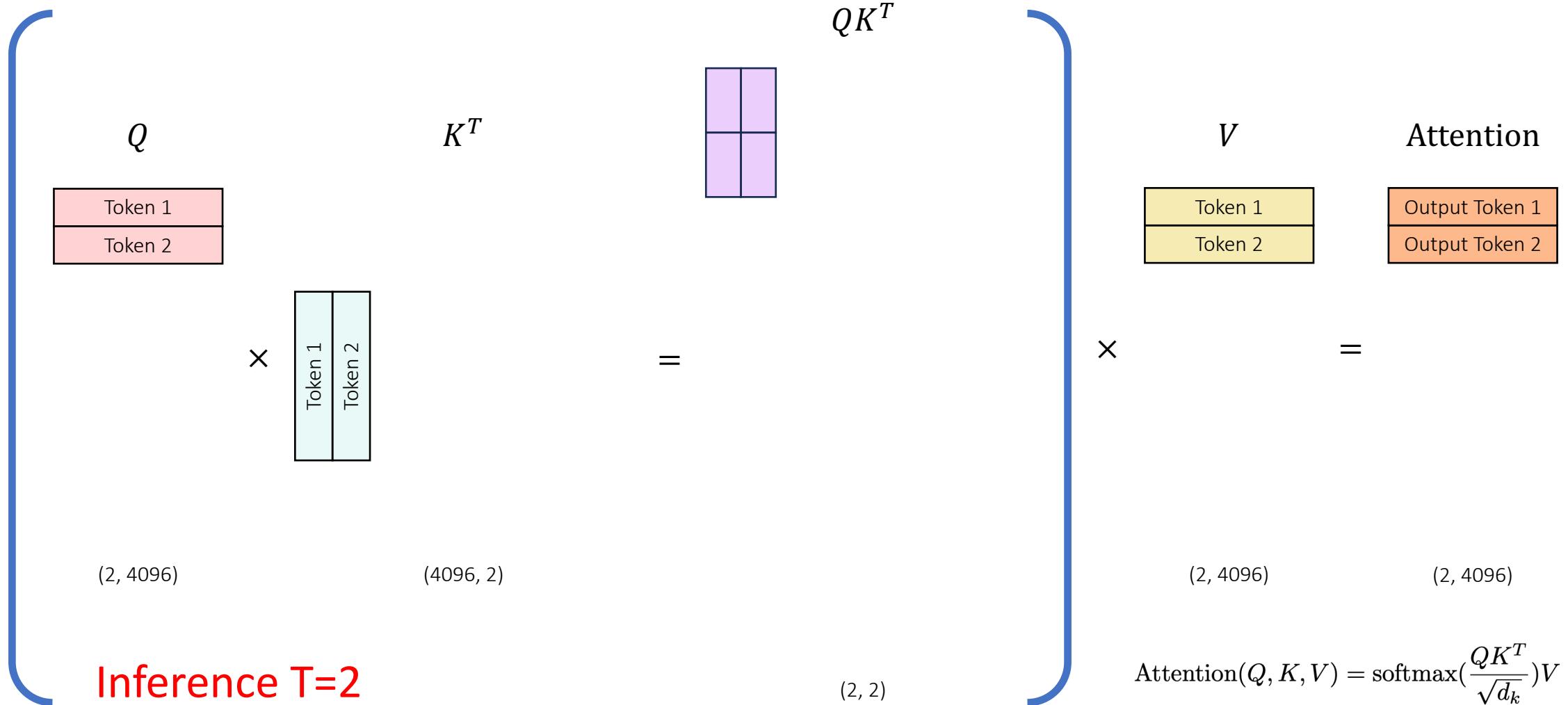
# Self-Attention during Next Token Prediction Task



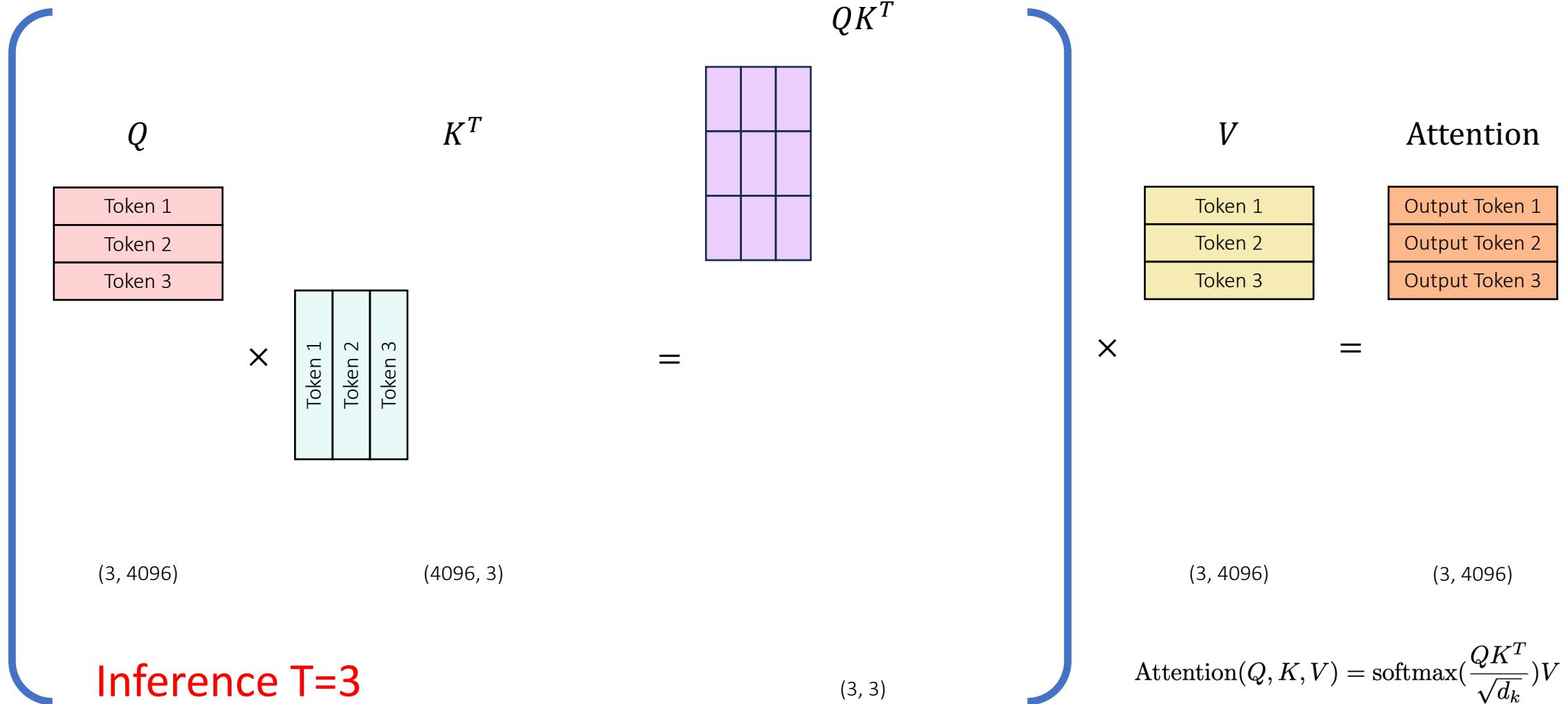
# Self-Attention during Next Token Prediction Task



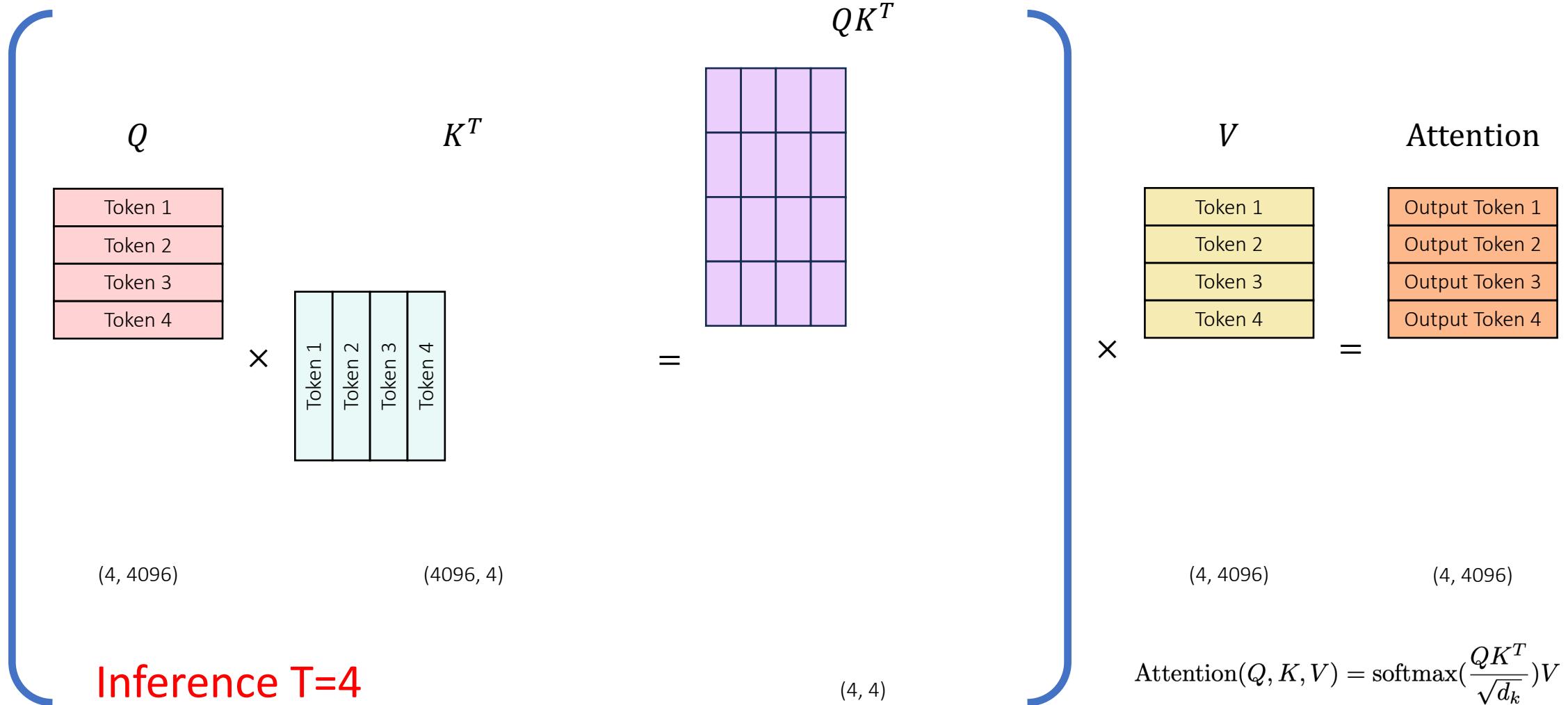
# Self-Attention during Next Token Prediction Task



# Self-Attention during Next Token Prediction Task



# Self-Attention during Next Token Prediction Task



1. We already computed these dot products in the previous steps. **Can we cache them?**

2. Since the model is causal, **we don't care about the attention of a token with its successors**, but only with the tokens before it.

3. **We don't care about these**, as we want to predict the next token and we have already predicted the previous ones.

$Q$

Token 1
Token 2
Token 3
Token 4

$\times$

$K^T$

Token 1
Token 2
Token 3
Token 4

(4, 4096)

(4096, 4)

Inference T=4

$QK^T$


$=$

4. We are only interested in this last row!

(4, 4)

$V$

Token 1
Token 2
Token 3
Token 4

$\times$

Attention

Output Token 1
Output Token 2
Output Token 3
Output Token 4

$=$

(4, 4096)

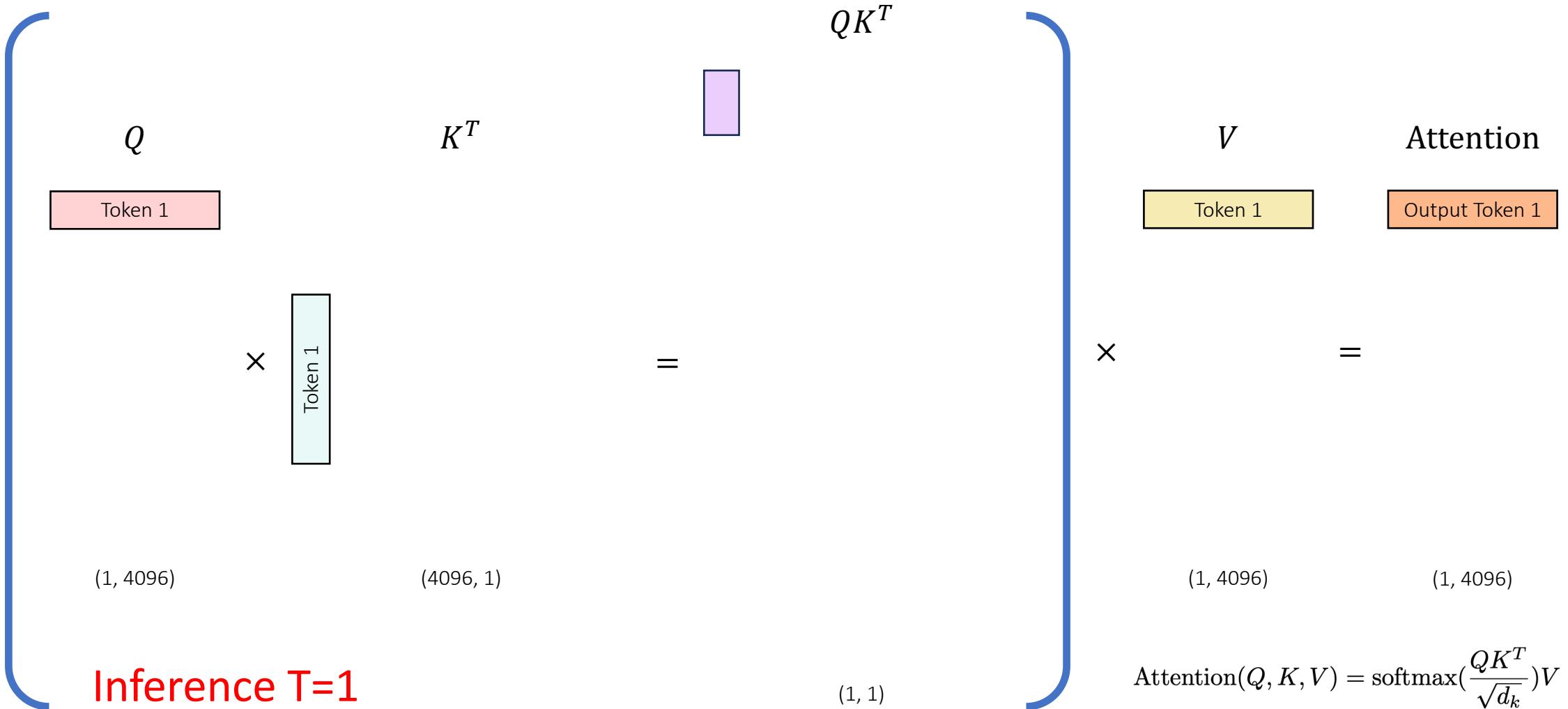
(4, 4096)

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

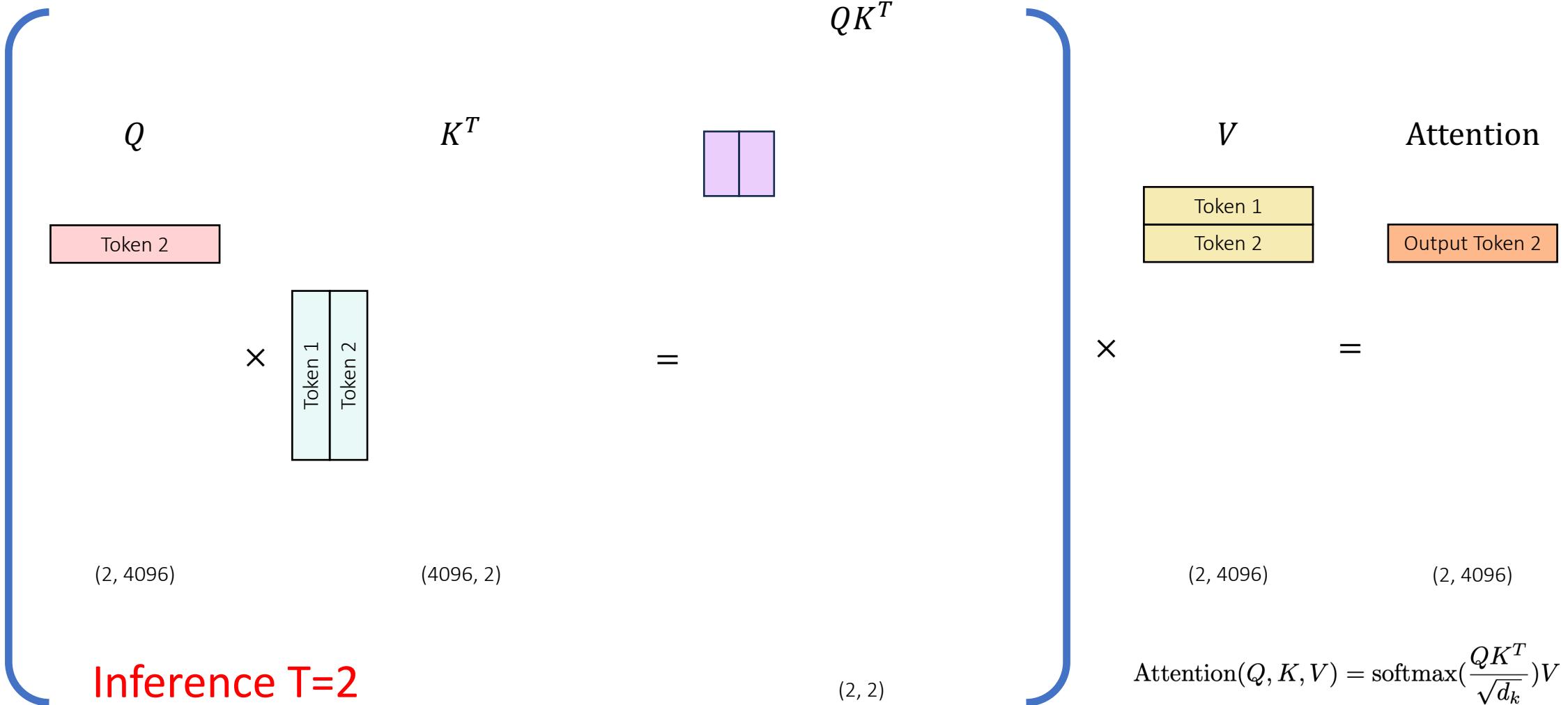
**ALL HAIL**

**THE KV CACHE**

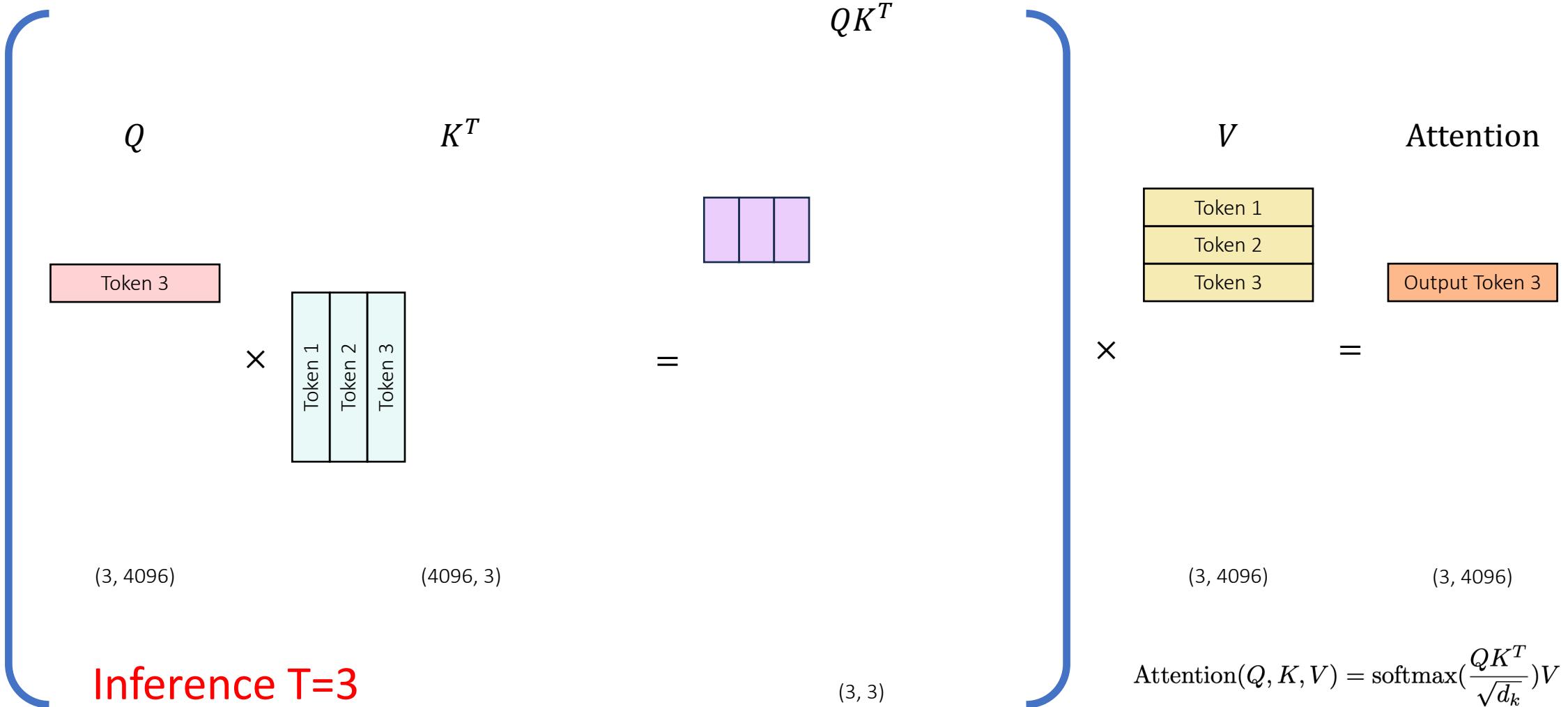
# Self-Attention with KV-Cache



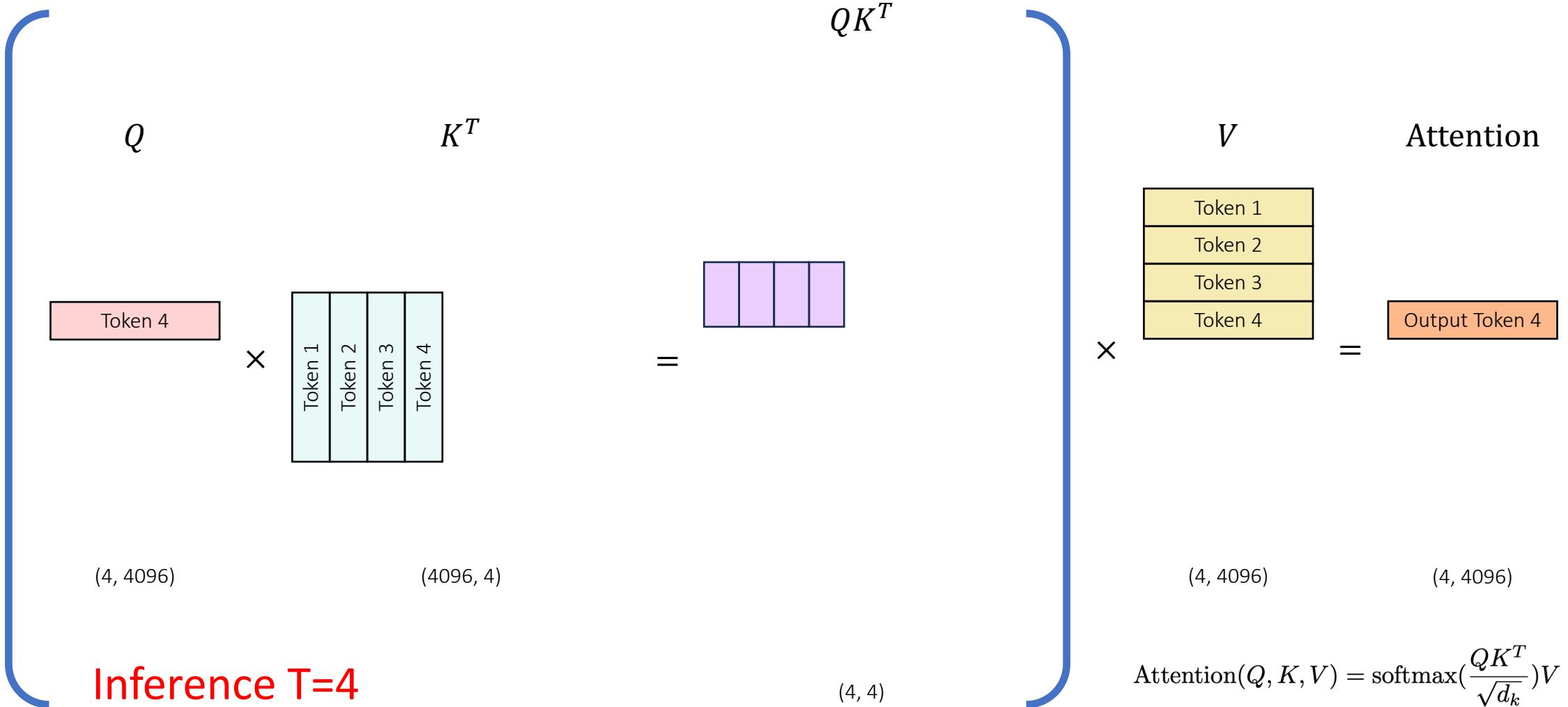
# Self-Attention with KV-Cache



# Self-Attention with KV-Cache

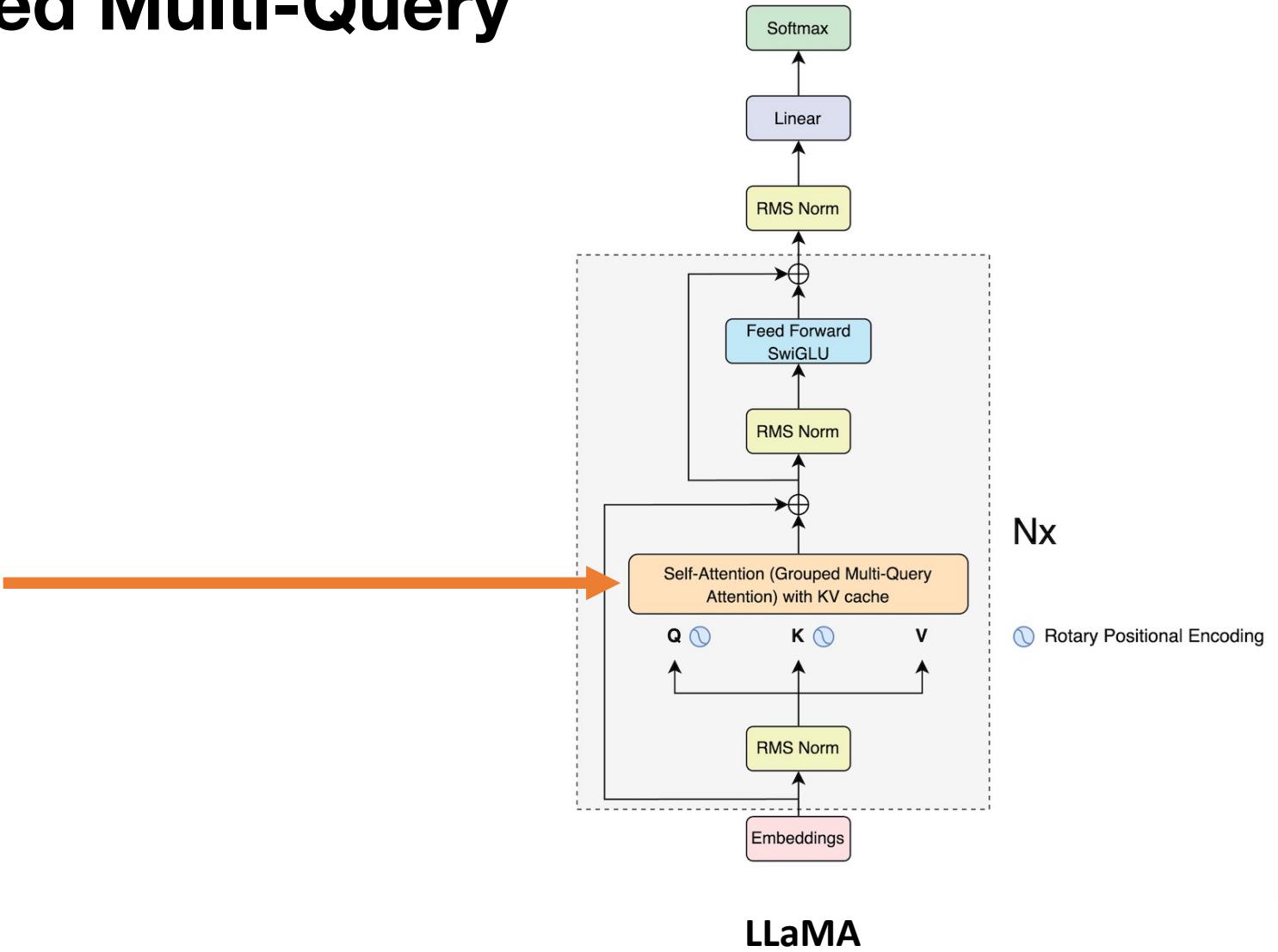


# Self-Attention with KV-Cache



# What is Grouped Multi-Query Attention?

Before we talk about Grouped Multi-Query Attention (MQA), we need to introduce its predecessor, the Multi-Query Attention (MQA).



# GPUs have a “problem”: They are too fast!

- In recent years, GPUs have become very fast at performing calculations, insomuch that the speed of computation (FLOPs) is much higher than the memory bandwidth (GB/s) or speed of data transfer between memory areas. For example, NVIDIA A100 can perform 19.5 TFLOPs while having a memory bandwidth of 2TB/s.
- This means that sometimes the bottleneck is not how many operations we perform, but how much data transfer our operations need, and that depends on the size and the quantity of the tensors involved in our calculations.
- For example, computing the same operation on the same tensor N times may be faster than computing the same operation on N different tensors, even if they have the same size, this is because the GPU may need to move the tensors around.
- **This means that our goal should not only be to optimize the number of operations we do, but also minimize the memory access/transfers that we perform.**

NVIDIA A100 TENSOR CORE GPU SPECIFICATIONS (SXM4 AND PCIE FORM FACTORS)				
	A100 40GB PCIe	A100 80GB PCIe	A100 40GB SXM	A100 80GB SXM
FP64		9.7 TFLOPS		
FP64 Tensor Core		19.5 TFLOPS		
FP32		19.5 TFLOPS		
Tensor Float 32 (TF32)		156 TFLOPS   312 TFLOPS*		
BFLOAT16 Tensor Core		312 TFLOPS   624 TFLOPS*		
FP16 Tensor Core		312 TFLOPS   624 TFLOPS*		
INT8 Tensor Core		624 TOPS   1248 TOPS*		
GPU Memory	40GB HBM2	80GB HBM2	40GB HBM2	80GB HBM2e
GPU Memory Bandwidth	1,555GB/s	1,935GB/s	1,555GB/s	2,039GB/s
Max Thermal Design Power (TDP)	250W	300W	400W	400W
Multi-Instance GPU	Up to 7 MIGs @ 5GB	Up to 7 MIGs @ 10GB	Up to 7 MIGs @ 5GB	Up to 7 MIGs @ 10GB
Form Factor	PCIe		SXM	
Interconnect	NVIDIA® NVLink® Bridge for 2 GPUs: 600GB/s ** PCIe Gen4: 64GB/s		NVLink: 600GB/s PCIe Gen4: 64GB/s	
Server Options	Partner and NVIDIA-Certified Systems™ with 1-8 GPUs		NVIDIA HGX™ A100-Partner and NVIDIA-Certified Systems with 4,8, or 16 GPUs NVIDIA DGX™ A100 with 8 GPUs	

\* With sparsity

\*\* SXM4 GPUs via HGX A100 server boards; PCIe GPUs via NVLink Bridge for up to two GPUs

# Introducing Multi-Query Attention

---

## Fast Transformer Decoding: One Write-Head is All You Need

---

Noam Shazeer  
Google  
[noam@google.com](mailto:noam@google.com)

November 7, 2019

# Comparing different attention algorithms: Vanilla Batched Multi-Head Attention

- Multi-head Attention as presented in the original paper “Attention is All You Need”
- By setting  $m = n$  (sequence length of query = sequence length of keys and values)
- The number of arithmetic operations performed is  $O(bnd^2)$
- The total memory involved in the operations, given by the sum of all the tensors involved in the calculations (including the derived ones!) is  $O(bnd + bhn^2 + d^2)$
- The ratio between the total memory and the number of arithmetic operations is  $O\left(\frac{1}{k} + \frac{1}{bn}\right)$
- In this case, the ratio is much smaller than 1, which means that the number of memory access we are performing is much less than the number of arithmetic operations, so the memory access is **not** the bottleneck here.

```
def MultiheadAttentionBatched():
    d, m, n, b, h, k, v = 512, 10, 10, 32, 8, (512 // 8), (512 // 8)

    X = torch.rand(b, n, d) # Query
    M = torch.rand(b, m, d) # Key and Value
    mask = torch.rand(b, h, n, m)
    P_q = torch.rand(h, d, k) # W_q
    P_k = torch.rand(h, d, k) # W_k
    P_v = torch.rand(h, d, v) # W_v
    P_o = torch.rand(h, d, v) # W_o

    Q = torch.einsum("bnd,hdk->bhnk ", X, P_q)
    K = torch.einsum("bmd,hdk->bhmk", M, P_k)
    V = torch.einsum("bmd,hdv->bhmv", M, P_v)

    logits = torch.einsum("bhnk,bhmk->bhnm", Q, K)
    weights = torch.softmax(logits + mask, dim=-1)

    O = torch.einsum("bhnm,bhmv->bhnv ", weights, V)
    Y = torch.einsum("bhnv,hdv->bnd ", O, P_o)

    return Y
```

# Comparing different attention algorithms: Batched Multi-Head Attention with KV-cache

- Uses the KV-cache to reduce the number of operations performed
- By setting  $m = n$  (sequence length of query = sequence length of keys and values)
- The number of arithmetic operations performed is  $O(bnd^2)$
- The total memory involved in the operations, given by the sum of all the tensors involved in the calculations (including the derived ones!) is  $O(bn^2d + nd^2)$
- The ratio between the total memory and the number of arithmetic operations is  $O\left(\frac{n}{d} + \frac{1}{b}\right)$
- When  $n \approx d$  (the sequence length is close to the size of the embedding vector) or  $b \approx 1$  (batch size is 1), the ratio becomes 1 and the memory access now becomes the bottleneck of the algorithm. For the batch size is not a problem, since it is generally much larger than 1, while for the  $\frac{n}{d}$  term, we need to reduce the sequence length.  
**But there's a better way...**

```
def MultiheadSelfAttentionIncremental():
    b, h, k, v = 512, 32, 8, (512 // 8), (512 // 8)

    m = 5 # Suppose we have already cached "m" tokens
    prev_K = torch.rand(b, h, m, k)
    prev_V = torch.rand(b, h, m, v)

    X = torch.rand(b, d) # Query
    M = torch.rand(b, d) # Key and Value
    P_q = torch.rand(h, d, k) # W_q
    P_k = torch.rand(h, d, k) # W_k
    P_v = torch.rand(h, d, v) # W_v
    P_o = torch.rand(h, d, v) # W_o

    q = torch.einsum("bd,hdk->bhk", X, P_q)
    new_K = torch.concat(
        [prev_K, torch.einsum("bd,hdk->bhk", M, P_k).unsqueeze(2)], axis=2
    )
    new_V = torch.concat(
        [prev_V, torch.einsum("bd,hdv->bhv", M, P_v).unsqueeze(2)], axis=2
    )
    logits = torch.einsum("bkh,bhmk->bhm", q, new_K)
    weights = torch.softmax(logits, dim=-1)
    O = torch.einsum("bhm,bhmv->bhv", weights, new_V)
    y = torch.einsum("bhv,hdv->bd", O, P_o)

    return y, new_K, new_V
```

# Comparing different attention algorithms: Multi-Query Attention with KV-cache

- We remove the  $h$  dimension from the  $K$  and the  $V$ , while keeping it for the  $Q$ . This means that all the different query heads will share the same keys and values.
- The number of arithmetic operations performed is  $O(bnd^2)$ .
- The total memory involved in the operations, given by the sum of all the tensors involved in the calculations (including the derived ones!) is  $O(bnd + bn^2k + nd^2)$ .
- The ratio between the total memory and the number of arithmetic operations is  $O\left(\frac{1}{d} + \frac{n}{nd} + \frac{1}{b}\right)$ .
- Comparing with the previous approach, we have reduced the expensive term  $\frac{n}{d}$  by a factor of  $h$ .
- The performance gains are important, while the model's quality degrades only a little bit.

```
def MultiquerySelfAttentionIncremental():
    d, b, h, k, v = 512, 32, 8, (512 // 8), (512 // 8)

    m = 5 # Suppose we have already cached "m" tokens
    prev_K = torch.rand(b, m, k)
    prev_V = torch.rand(b, m, v)

    X = torch.rand(b, d) # Query
    M = torch.rand(b, d) # Key and Value
    P_q = torch.rand(h, d, k) # W_q
    P_k = torch.rand(d, k) # W_k
    P_v = torch.rand(d, v) # W_v
    P_o = torch.rand(h, d, v) # W_o

    q = torch.einsum("bd,hdk->bhk", X, P_q)
    K = torch.concat([prev_K, torch.einsum("bd,dk->bk", M, P_k).unsqueeze(1)], axis=1)
    V = torch.concat([prev_V, torch.einsum("bd,dv->bv", M, P_v).unsqueeze(1)], axis=1)
    logits = torch.einsum("bhk,bmk->bhm", q, K)
    weights = torch.softmax(logits, dim=-1)
    O = torch.einsum("bhm,bmv->bhv", weights, V)
    y = torch.einsum("bhv,hdv->bd", O, P_o)

    return y, K, V
```

# Speed and quality comparisons

BLEU score on a translation task (English - German)

Table 1: WMT14 EN-DE Results.

Attention Type	$h$	$d_k, d_v$	$d_{ff}$	ln(PPL) (dev)	BLEU (dev)	BLEU (test) beam 1 / 4
multi-head	8	128	4096	<b>1.424</b>	<b>26.7</b>	27.7 / 28.4
multi-query	8	128	5440	1.439	26.5	27.5 / <b>28.5</b>
multi-head local	8	128	4096	1.427	26.6	27.5 / 28.3
multi-query local	8	128	5440	1.437	26.5	27.6 / 28.2
multi-head	1	128	6784	1.518	25.8	
multi-head	2	64	6784	1.480	26.2	26.8 / 27.9
multi-head	4	32	6784	1.488	26.1	
multi-head	8	16	6784	1.513	25.8	

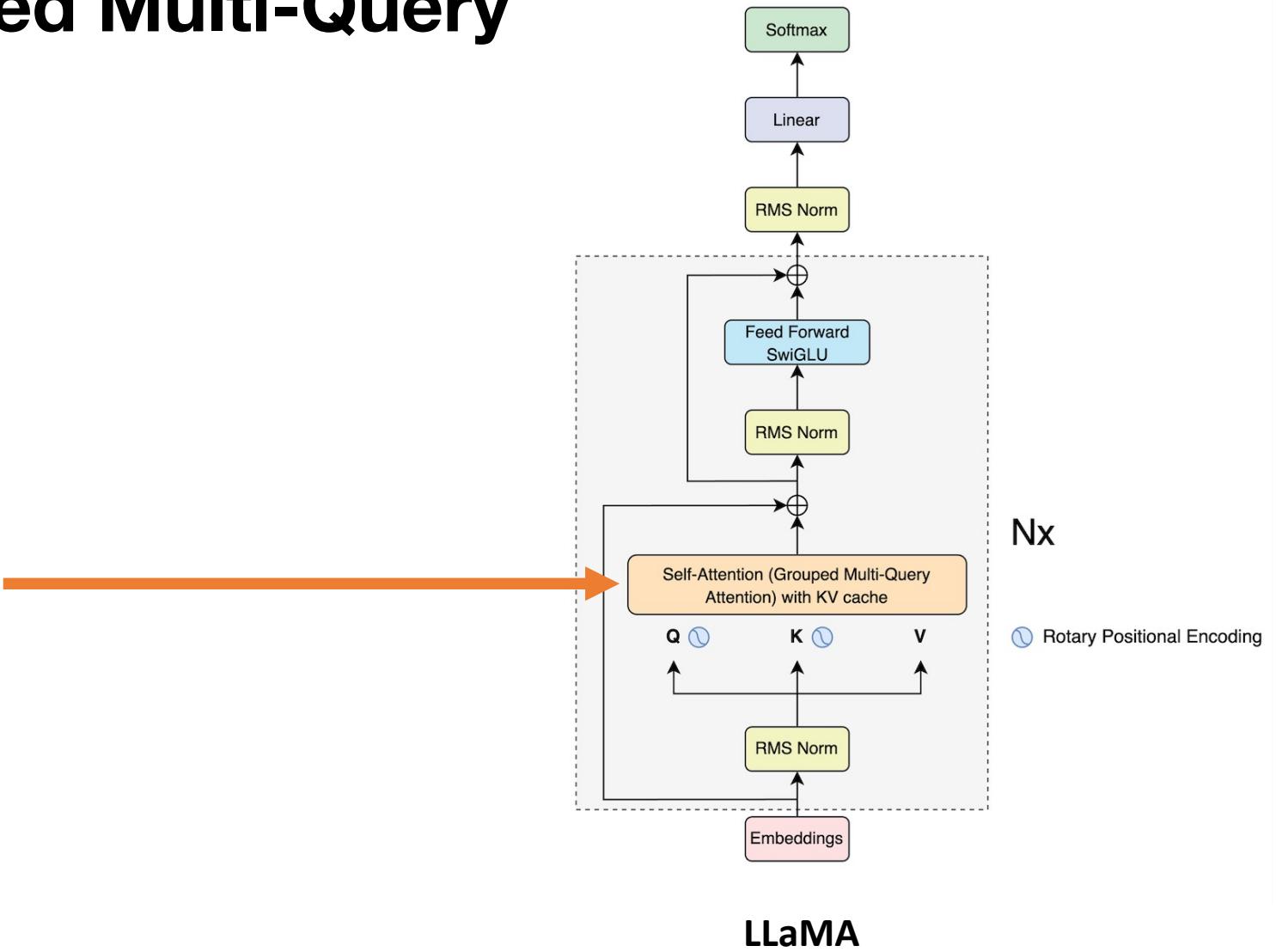
Table 2: Amortized training and inference costs for WMT14 EN-DE Translation Task with sequence length 128. Values listed are in TPUv2-microseconds per output token.

Attention Type	Training	Inference enc. + dec.	Beam-4 Search enc. + dec.
multi-head	13.2	1.7 + 46	2.0 + 203
multi-query	<b>13.0</b>	1.5 + 3.8	1.6 + 32
multi-head local	13.2	1.7 + 23	1.9 + 47
multi-query local	<b>13.0</b>	<b>1.5 + 3.3</b>	<b>1.6 + 16</b>

To demonstrate that local-attention and multi-query attention are orthogonal, we also trained "local" versions of the baseline and multi-query models, where the decoder-self-attention layers (but not the other attention layers) restrict attention to the current position and the previous 31 positions.

# What is Grouped Multi-Query Attention?

Now, let's talk about Grouped MQA!



# Grouped Multi-Query Attention: A compromise between two extremes

## Multi-Head Attention

- High quality
- Computationally slow

## Grouped Multi-Query Attention

- A good compromise between quality and speed

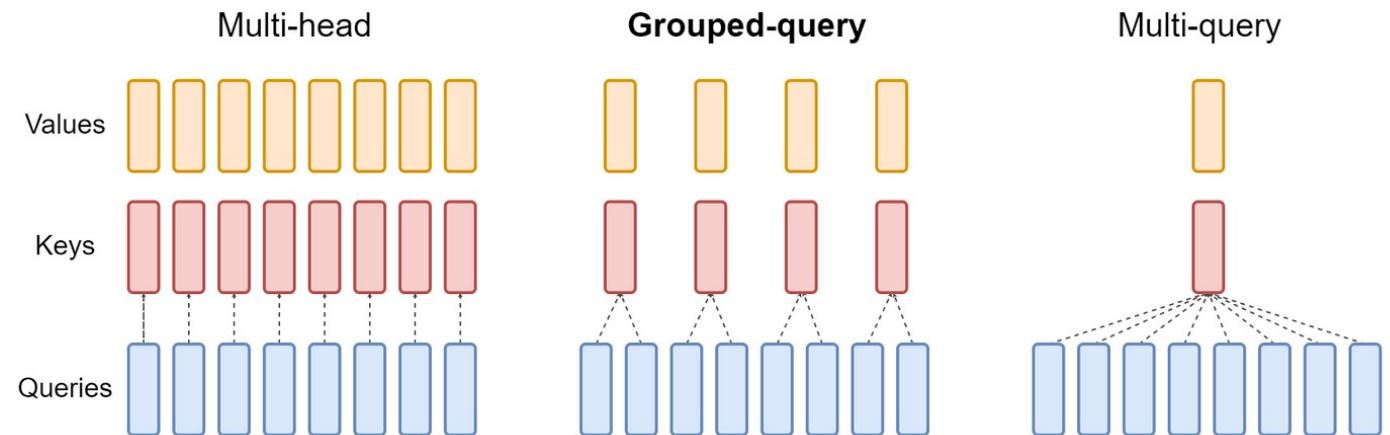
## Multi-Query Attention

- Loss in quality
- Computationally fast

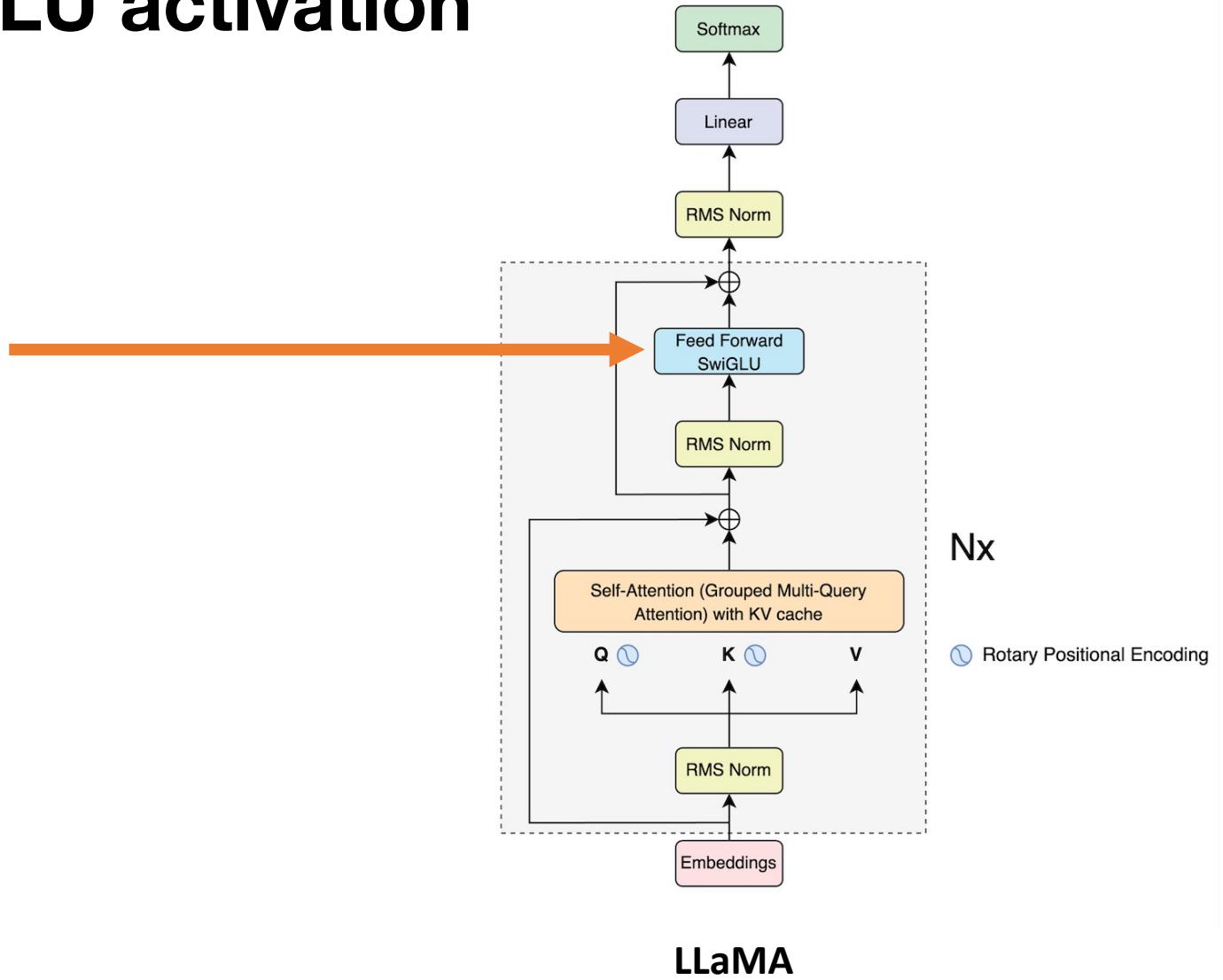
GQA: Training Generalized Multi-Query Transformer Models from Multi-Head Checkpoints

Joshua Ainslie\*, James Lee-Thorp\*, Michiel de Jong\*<sup>†</sup>  
Yury Zemlyanskiy, Federico Lebrón, Sumit Sanghi

Google Research



# What is the SwiGLU activation function?



# SwiGLU Activation Function

---

## GLU Variants Improve Transformer

---

Noam Shazeer  
Google  
[noam@google.com](mailto:noam@google.com)

February 14, 2020

# SwiGLU Activation Function

- The author compared the performance of a Transformer model by using different activation functions in the Feed-Forward layer of the Transformer architecture.

$$\begin{aligned} \text{ReLU}(x, W, V, b, c) &= \max(0, xW + b) \otimes (xV + c) \\ \text{GEGLU}(x, W, V, b, c) &= \text{GELU}(xW + b) \otimes (xV + c) \\ \text{SwiGLU}(x, W, V, b, c, \beta) &= \text{Swish}_\beta(xW + b) \otimes (xV + c) \end{aligned} \tag{5}$$

In this paper, we propose additional variations on the Transformer FFN layer which use GLU or one of its variants in place of the first linear transformation and the activation function. Again, we omit the bias terms.

$$\begin{aligned} \text{FFN}_{\text{GLU}}(x, W, V, W_2) &= (\sigma(xW) \otimes xV)W_2 \\ \text{FFN}_{\text{Bilinear}}(x, W, V, W_2) &= (xW \otimes xV)W_2 \\ \text{FFN}_{\text{ReLU}}(x, W, V, W_2) &= (\max(0, xW) \otimes xV)W_2 \\ \text{FFN}_{\text{GEGLU}}(x, W, V, W_2) &= (\text{GELU}(xW) \otimes xV)W_2 \\ \text{FFN}_{\text{SwiGLU}}(x, W, V, W_2) &= (\text{Swish}_1(xW) \otimes xV)W_2 \end{aligned} \tag{6}$$

All of these layers have three weight matrices, as opposed to two for the original FFN. To keep the number of parameters and the amount of computation constant, we reduce the number of hidden units  $d_{ff}$  (the second dimension of  $W$  and  $V$  and the first dimension of  $W_2$ ) by a factor of  $\frac{2}{3}$  when comparing these layers to the original two-matrix version.

# SwiGLU Activation Function

**Transformer** ("Attention is all you need")

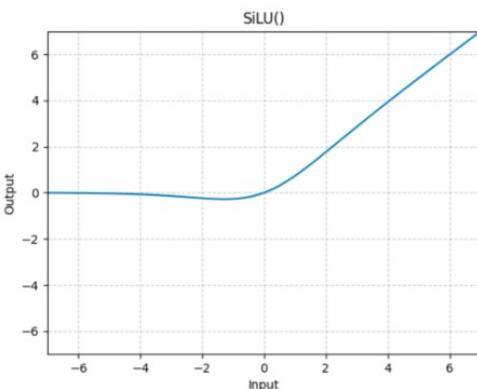
$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

**LLaMA**

$$\text{FFN}_{\text{SwiGLU}}(x, W, V, W_2) = (\text{Swish}_1(xW) \otimes xV)W_2$$

We use the swish function with  $\beta = 1$ . In this case it's called the **Sigmoid Linear Unit (SiLU)** function.

$$\text{swish}(x) = x \text{ sigmoid}(\beta x) = \frac{x}{1 + e^{-\beta x}}$$



```
class FeedForward(nn.Module):
    def __init__(self,
                 dim,
                 hidden_dim: int,
                 multiple_of: int,
                 ffn_dim_multiplier: Optional[float],
                 ):
        super().__init__()
        hidden_dim = int(2 * hidden_dim / 3)
        # custom dim factor multiplier
        if ffn_dim_multiplier is not None:
            hidden_dim = int(ffn_dim_multiplier * hidden_dim)
        hidden_dim = multiple_of * ((hidden_dim + multiple_of - 1) // multiple_of)

        self.w1 = ColumnParallelLinear(
            dim, hidden_dim, bias=False, gather_output=False, init_method=lambda x: x
        )
        self.w2 = RowParallelLinear(
            hidden_dim, dim, bias=False, input_is_parallel=True, init_method=lambda x: x
        )
        self.w3 = ColumnParallelLinear(
            dim, hidden_dim, bias=False, gather_output=False, init_method=lambda x: x
        )

    def forward(self, x):
        return self.w2(F.silu(self.w1(x)) * self.w3(x))
```

# How well does it perform?

Table 1: Heldout-set log-perplexity for Transformer models on the segment-filling task from [Raffel et al., 2019]. All models are matched for parameters and computation.

Training Steps	65,536	524,288
FFN <sub>ReLU</sub> (baseline)	1.997 (0.005)	1.677
FFN <sub>GELU</sub>	1.983 (0.005)	1.679
FFN <sub>Swish</sub>	1.994 (0.003)	1.683
FFN <sub>GLU</sub>	1.982 (0.006)	1.663
FFN <sub>Bilinear</sub>	1.960 (0.005)	1.648
FFN <sub>GEGLU</sub>	<b>1.942</b> (0.004)	<b>1.633</b>
FFN <sub>SwiGLU</sub>	<b>1.944</b> (0.010)	<b>1.636</b>
FFN <sub>ReGLU</sub>	1.953 (0.003)	1.645



Table 2: GLUE Language-Understanding Benchmark [Wang et al., 2018] (dev).

	Score Average	CoLA	SST-2	MRPC	MRPC	STSB	STSB	QQP	QQP	MNLI <sub>m</sub>	MNLI <sub>mm</sub>	QNLI	RTE
		MCC	Acc	F1	Acc	PCC	SCC	F1	Acc	Acc	Acc	Acc	Acc
FFN <sub>ReLU</sub>	83.80	51.32	94.04	<b>93.08</b>	<b>90.20</b>	89.64	89.42	89.01	91.75	85.83	86.42	92.81	80.14
FFN <sub>GELU</sub>	83.86	53.48	94.04	92.81	<b>90.20</b>	89.69	89.49	88.63	91.62	85.89	86.13	92.39	80.51
FFN <sub>Swish</sub>	83.60	49.79	93.69	92.31	89.46	89.20	88.98	88.84	91.67	85.22	85.02	92.33	81.23
FFN <sub>GLU</sub>	84.20	49.16	94.27	92.39	89.46	89.46	89.35	88.79	91.62	86.36	86.18	92.92	<b>84.12</b>
FFN <sub>GEGLU</sub>	84.12	53.65	93.92	92.68	89.71	90.26	90.13	89.11	91.85	86.15	86.17	92.81	79.42
FFN <sub>Bilinear</sub>	83.79	51.02	<b>94.38</b>	92.28	89.46	90.06	89.84	88.95	91.69	<b>86.90</b>	<b>87.08</b>	92.92	81.95
FFN <sub>SwiGLU</sub>	84.36	51.59	93.92	92.23	88.97	<b>90.32</b>	<b>90.13</b>	<b>89.14</b>	<b>91.87</b>	86.45	86.47	<b>92.93</b>	83.39
FFN <sub>ReGLU</sub>	<b>84.67</b>	<b>56.16</b>	<b>94.38</b>	92.06	89.22	89.97	89.85	88.86	91.72	86.20	86.40	92.68	81.59
[Raffel et al., 2019]	83.28	53.84	92.68	92.07	88.92	88.02	87.94	88.67	91.56	84.24	84.57	90.48	76.28
ibid. stddev.	0.235	1.111	0.569	0.729	1.019	0.374	0.418	0.108	0.070	0.291	0.231	0.361	1.393



# Why SwiGLU works so well?

## 4 Conclusions

We have extended the GLU family of layers and proposed their use in Transformer. In a transfer-learning setup, the new variants seem to produce better perplexities for the de-noising objective used in pre-training, as well as better results on many downstream language-understanding tasks. These architectures are simple to implement, and have no apparent computational drawbacks. We offer no explanation as to why these architectures seem to work; we attribute their success, as all else, to divine benevolence.

