# LLM NOTES

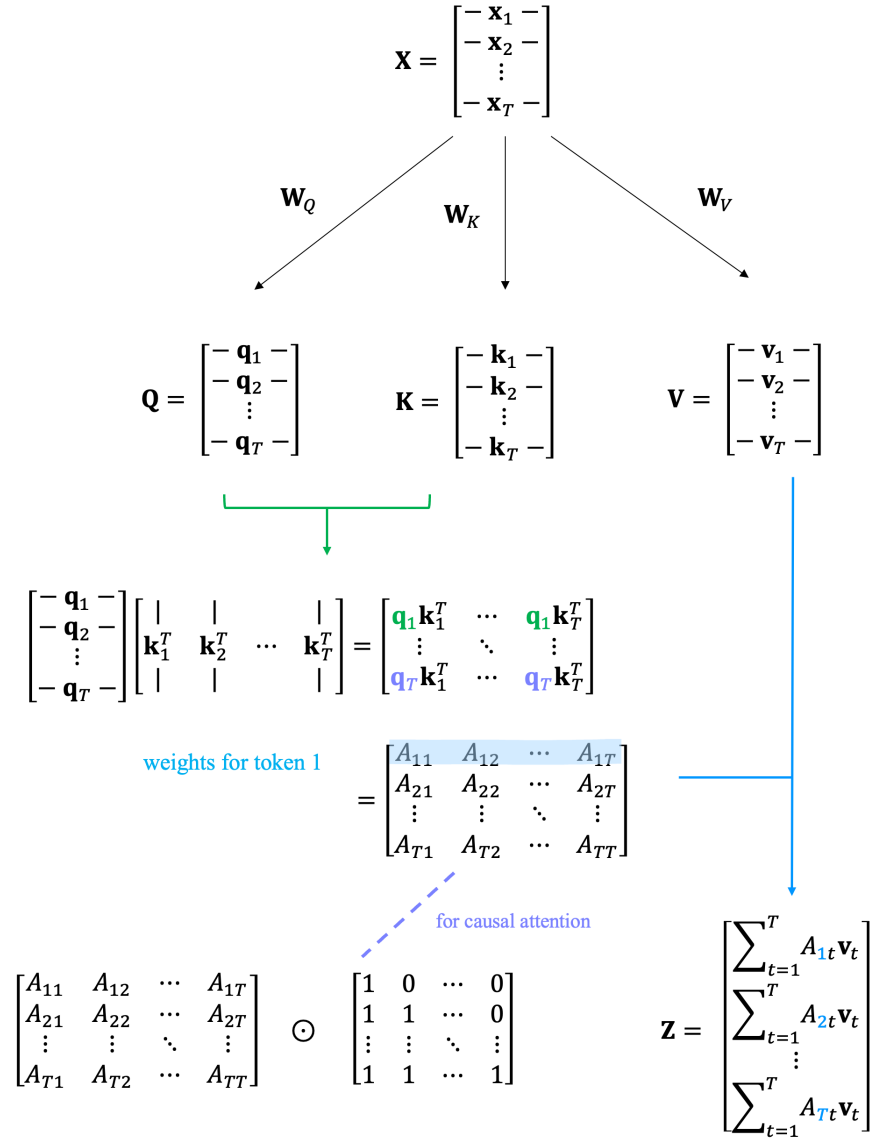### A PREPRINT

July 18, 2025

## Contents

# 1 Attention

## 1.1 Self Attention

Given input $\boldsymbol{X} \in \mathbb{R}^{T \times D}$ where $T$ is the sequence length and $D$ is the embedding dimension, we have

$$\boldsymbol{Q} = \boldsymbol{W}_Q \boldsymbol{X} \in \mathbb{R}^{T \times d_k} \quad \boldsymbol{K} = \boldsymbol{W}_K \boldsymbol{X} \in \mathbb{R}^{T \times d_k} \quad \boldsymbol{V} = \boldsymbol{W}_V \boldsymbol{X} \in \mathbb{R}^{T \times d_v} \tag{1}$$

$$\text{Attn}(\boldsymbol{Q}, \boldsymbol{K}, \boldsymbol{V}) = \text{softmax}\left(\frac{\boldsymbol{Q}\boldsymbol{K}^\top}{\sqrt{d_k}}\right)\boldsymbol{V} \tag{2}$$

$$\mathbf{X} = \begin{bmatrix} - \mathbf{x}_1 - \\ - \mathbf{x}_2 - \\ \vdots \\ - \mathbf{x}_T - \end{bmatrix}$$

$\mathbf{W}_Q \qquad \mathbf{W}_K \qquad \mathbf{W}_V$

$$\mathbf{Q} = \begin{bmatrix} - \mathbf{q}_1 - \\ - \mathbf{q}_2 - \\ \vdots \\ - \mathbf{q}_T - \end{bmatrix} \qquad \mathbf{K} = \begin{bmatrix} - \mathbf{k}_1 - \\ - \mathbf{k}_2 - \\ \vdots \\ - \mathbf{k}_T - \end{bmatrix} \qquad \mathbf{V} = \begin{bmatrix} - \mathbf{v}_1 - \\ - \mathbf{v}_2 - \\ \vdots \\ - \mathbf{v}_T - \end{bmatrix}$$

$$\begin{bmatrix} - \mathbf{q}_1 - \\ - \mathbf{q}_2 - \\ \vdots \\ - \mathbf{q}_T - \end{bmatrix} \begin{bmatrix} | & | & & | \\ \mathbf{k}_1^T & \mathbf{k}_2^T & \cdots & \mathbf{k}_T^T \\ | & | & & | \end{bmatrix} = \begin{bmatrix} \mathbf{q}_1\mathbf{k}_1^T & \cdots & \mathbf{q}_1\mathbf{k}_T^T \\ \vdots & \ddots & \vdots \\ \mathbf{q}_T\mathbf{k}_1^T & \cdots & \mathbf{q}_T\mathbf{k}_T^T \end{bmatrix}$$

weights for token 1

$$= \begin{bmatrix} A_{11} & A_{12} & \cdots & A_{1T} \\ A_{21} & A_{22} & \cdots & A_{2T} \\ \vdots & \vdots & \ddots & \vdots \\ A_{T1} & A_{T2} & \cdots & A_{TT} \end{bmatrix}$$

for causal attention

$$\begin{bmatrix} A_{11} & A_{12} & \cdots & A_{1T} \\ A_{21} & A_{22} & \cdots & A_{2T} \\ \vdots & \vdots & \ddots & \vdots \\ A_{T1} & A_{T2} & \cdots & A_{TT} \end{bmatrix} \odot \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 1 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \cdots & 1 \end{bmatrix} \qquad \mathbf{Z} = \begin{bmatrix} \sum_{t=1}^{T} A_{1t}\mathbf{v}_t \\ \sum_{t=1}^{T} A_{2t}\mathbf{v}_t \\ \vdots \\ \sum_{t=1}^{T} A_{Tt}\mathbf{v}_t \end{bmatrix}$$
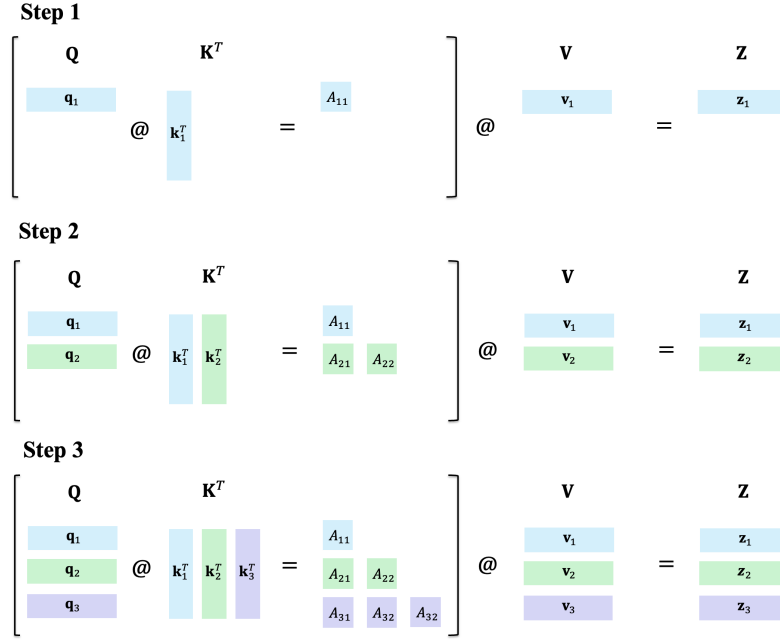
## 1.2 Multi-Head Attention

In multi-head attention, input $X$ is first passed through $H$ self-attention layer in parallel. Then, the output from each head is concatenated together and fused by a linear projection

$$\left[Z^{(1)}, Z^{(2)}, \ldots, Z^{(H)}\right] W^O = \begin{bmatrix} z_1^{(1)} & z_1^{(2)} & \ldots & z_1^{(H)} \\ z_2^{(1)} & z_2^{(2)} & \ldots & z_2^{(H)} \\ \vdots & \vdots & \ldots & \vdots \\ z_T^{(1)} & z_T^{(2)} & \ldots & z_T^{(H)} \end{bmatrix} W^O \tag{3}$$
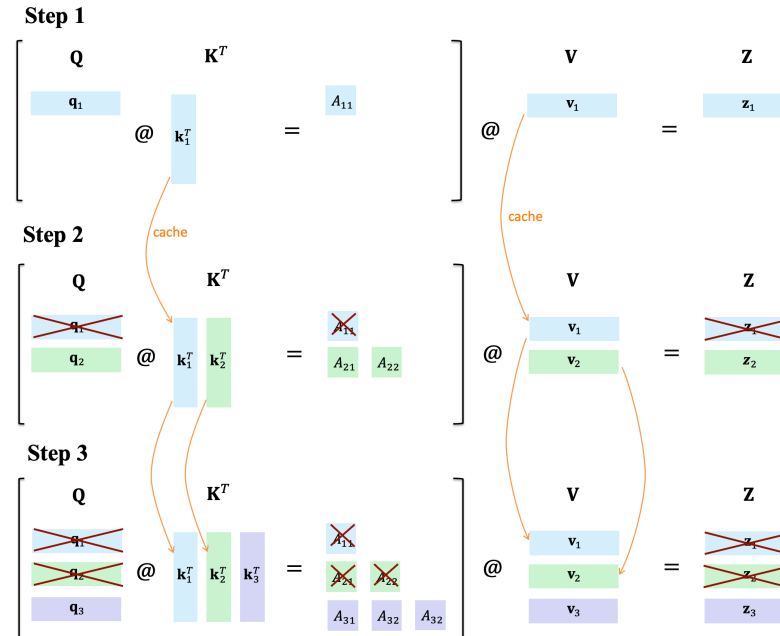
## 1.3  KV-Cache

During inference we still use causal masking because this is how the model being trained. Let's look at a simple case where we only give the model a start token <s> and asks it to generate stuff:

**Step 1**

$$\mathbf{Q} \quad \mathbf{K}^T \qquad = \qquad A_{11} \qquad @ \qquad \mathbf{V} \qquad = \qquad \mathbf{Z}$$

$q_1 \quad @ \quad k_1^T \qquad A_{11} \qquad @ \qquad v_1 \qquad = \qquad z_1$

**Step 2**

$q_1$
$q_2 \quad @ \quad k_1^T \; k_2^T \qquad A_{11} / A_{21} \; A_{22} \qquad @ \qquad v_1 / v_2 \qquad = \qquad z_1 / z_2$

**Step 3**

$q_1$
$q_2 \quad @ \quad k_1^T \; k_2^T \; k_3^T \qquad A_{11} / A_{21} A_{22} / A_{31} A_{32} A_{32} \qquad @ \qquad v_1 / v_2 / v_3 \qquad = \qquad z_1 / z_2 / z_3$
$q_3$

KV-cache is built on this two observations:

- At each time step $t$, due to causal masking $\boldsymbol{k}_{<t}$ and $\boldsymbol{v}_{<t}$ will remain the same
- To predict <token$_{t+1}$> we only need embedding of <token$_t$>.

Therefore, we can make prediction efficiently by drop redundant and unnecessary computation

**Step 1**

$\mathbf{Q} \qquad \mathbf{K}^T \qquad = \qquad A_{11} \qquad @ \qquad \mathbf{V} \qquad = \qquad \mathbf{Z}$

$q_1 \quad @ \quad k_1^T \qquad A_{11} \qquad @ \qquad v_1 \qquad z_1$

*cache*                    *cache*

**Step 2**

$\mathbf{Q} \qquad \mathbf{K}^T \qquad = \qquad @ \qquad \mathbf{V} \qquad = \qquad \mathbf{Z}$

$q_2 \quad @ \quad k_1^T \; k_2^T \qquad A_{21} \; A_{22} \qquad @ \qquad v_2 \qquad z_2$

**Step 3**

$\mathbf{Q} \qquad \mathbf{K}^T \qquad = \qquad @ \qquad \mathbf{V} \qquad = \qquad \mathbf{Z}$

$q_3 \quad @ \quad k_1^T \; k_2^T \; k_3^T \qquad A_{31} \; A_{32} \; A_{32} \qquad @ \qquad v_3 \qquad z_3$

Basically we have

Token 1: [K1, V1] $\rightarrow$ Cache: [K1, V1]

Token 2: [K2, V2] $\rightarrow$ Cache: [K1, K2], [V1, V2]

...

Token n: [Kn, Vn] $\rightarrow$ Cache: [K1, K2, ..., Kn], [V1, V2, ..., Vn]

## 2 Positional Embedding

# References