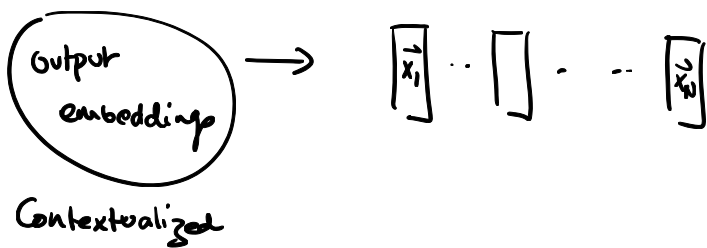
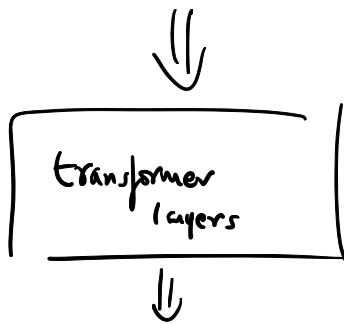
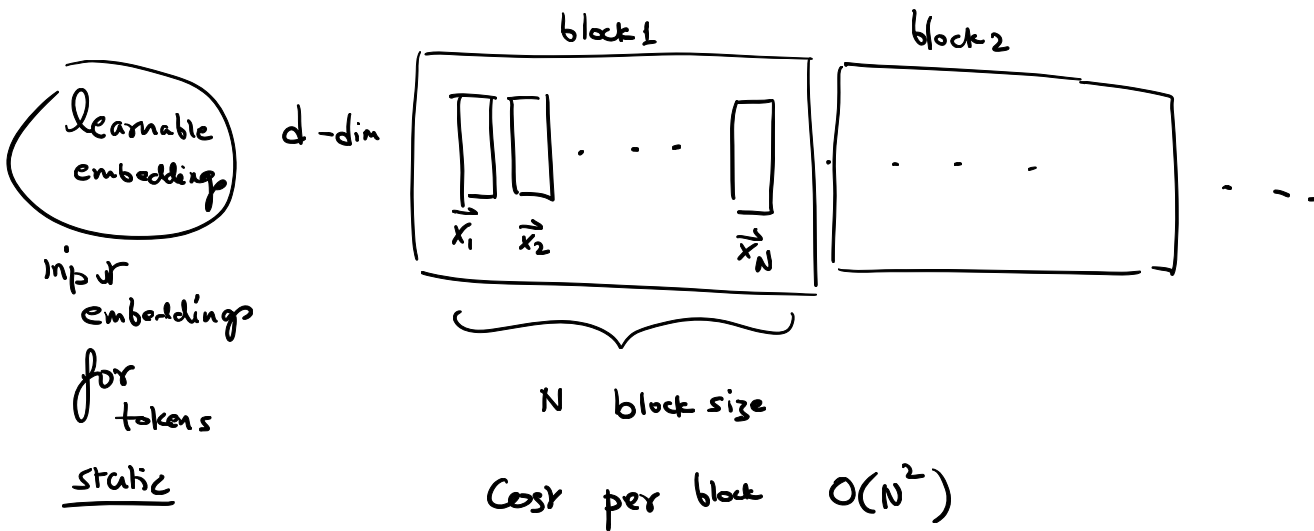


Transformers

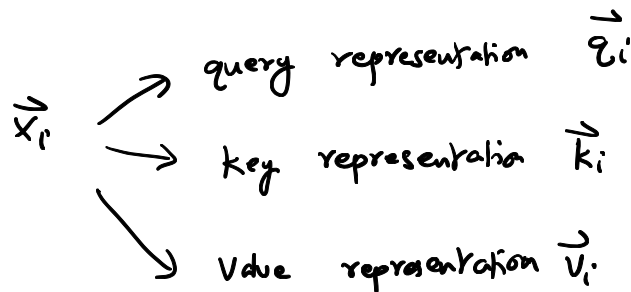
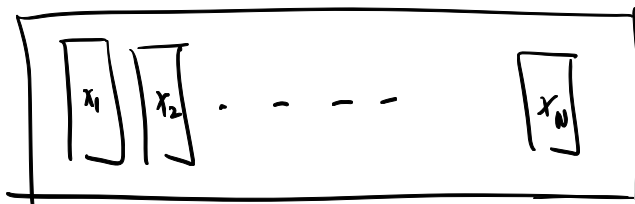
dot-product attention

Input: block / "sequence"



Same token
different embeddings
in diff blocks

Self-attention



Dimensionality

$$\vec{x}_i : d_{\text{model}} (512)$$

$$\vec{q}_i : d_k (512/p = 64)$$

$$\vec{k}_i : d_k (64)$$

$$\vec{v}_i : d_v (64)$$

$\vec{x}_i \leftarrow$ attention

\vec{q}_i vs all keys \vec{k}_j for all $j=1, \dots, N$

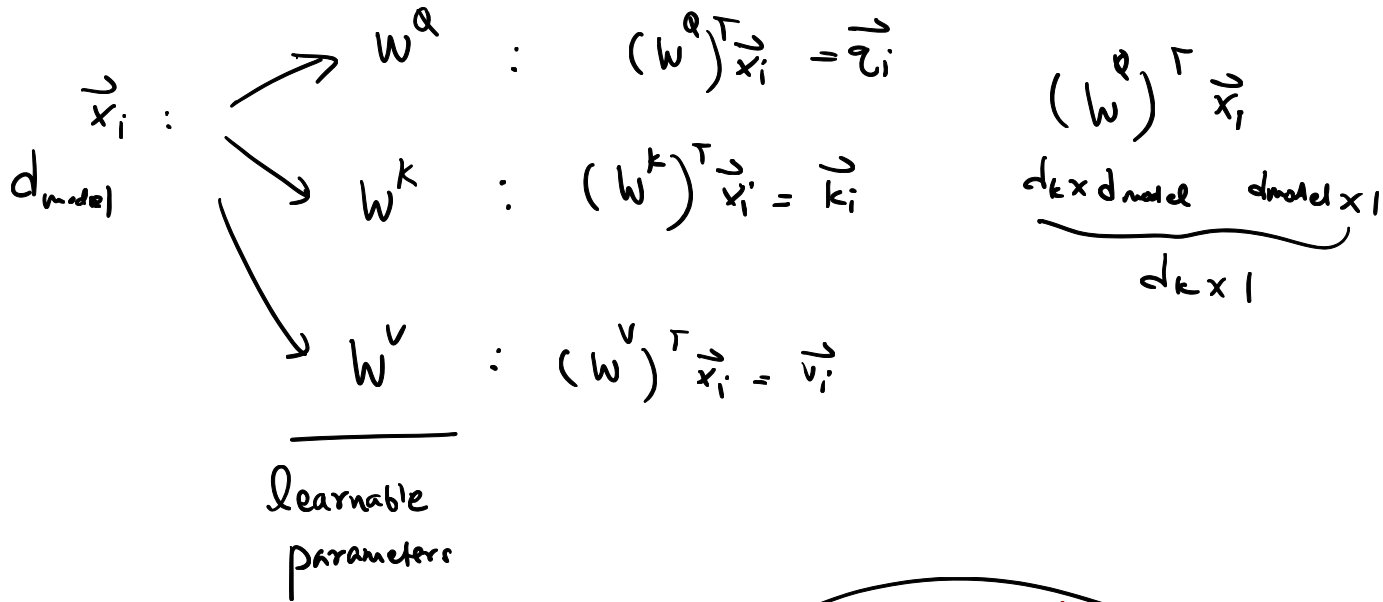
$$\text{attention vector} : \alpha_i = \text{Softmax} \left(\frac{\vec{q}_i^T \vec{k}_1}{\sqrt{d_k}}, \frac{\vec{q}_i^T \vec{k}_2}{\sqrt{d_k}}, \dots, \frac{\vec{q}_i^T \vec{k}_N}{\sqrt{d_k}} \right)$$

$\sqrt{d_k}$: rescaling of dot product

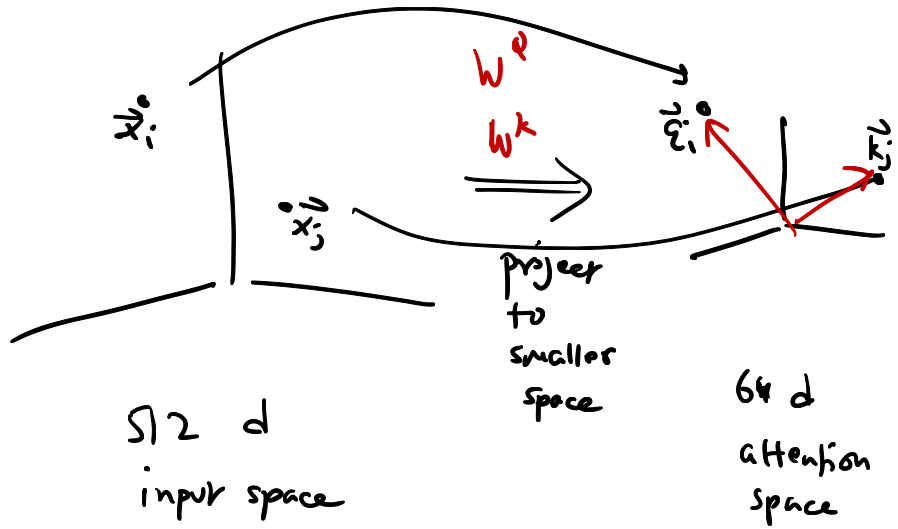
α_{ij} : attention \vec{x}_j pays to \vec{x}_i
(\vec{k}_j) (\vec{q}_i)

Update the value representation \vec{v}_i : weighted sum of all value vectors (in the block)

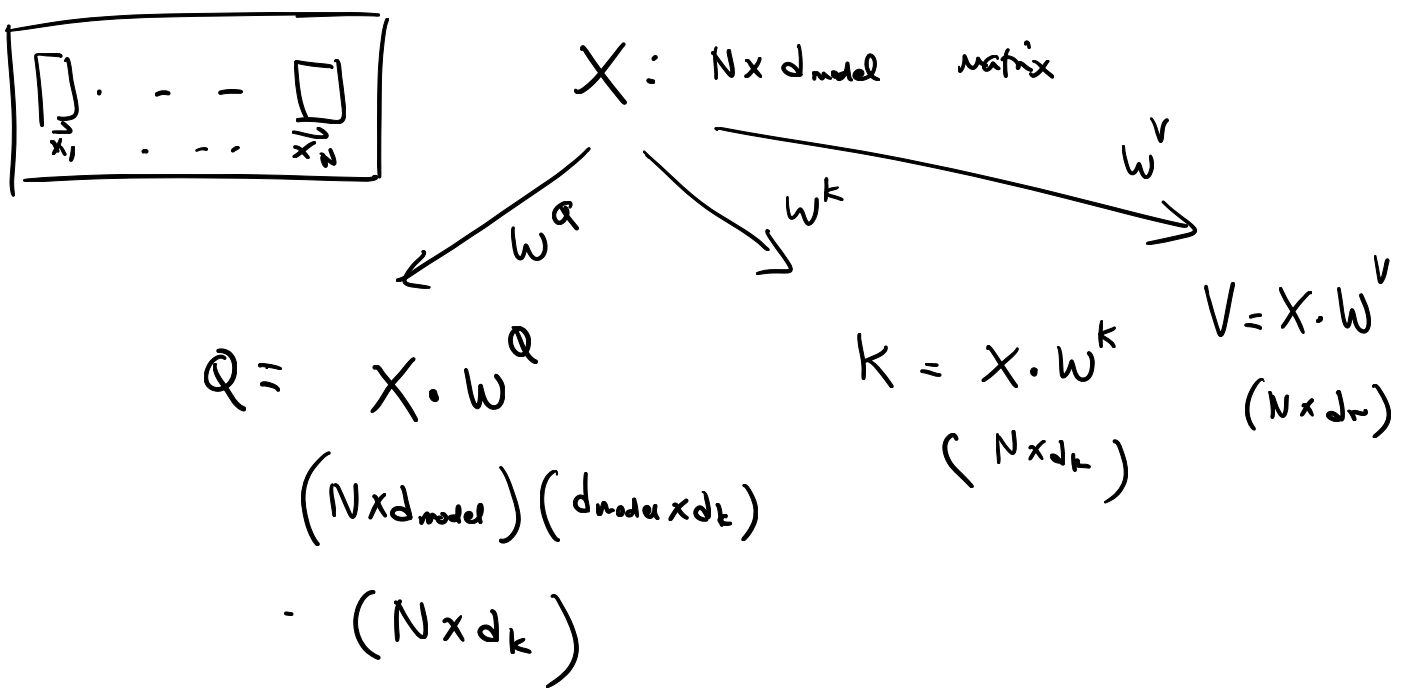
$$\vec{v}_i = \sum_{j=1}^N \alpha_{ij} \vec{v}_j$$

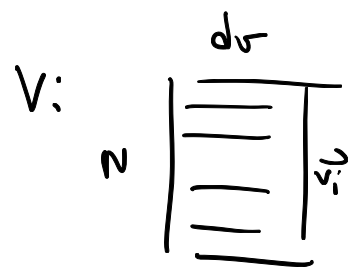
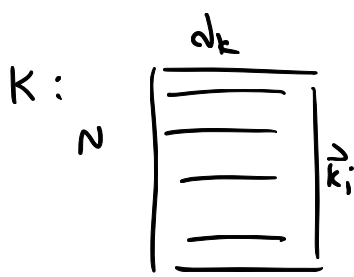
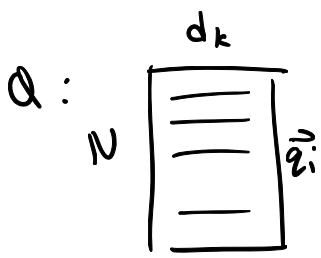


w^q : $d_{model} \times d_k$
 w^k : $d_{model} \times d_k$
 w^v : $d_{model} \times d_v$

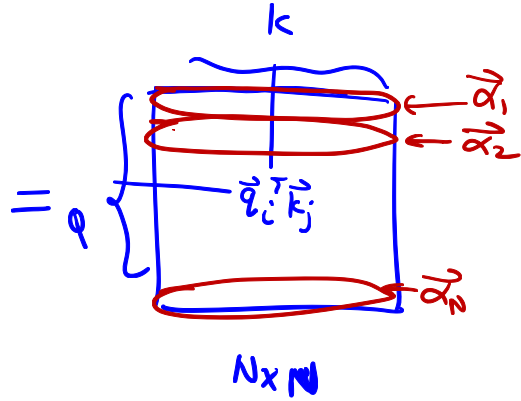
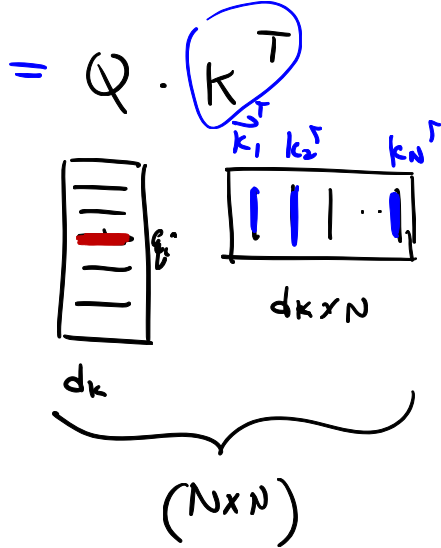


Matrix Approach : All attention for all elements of the block





all pair-wise dot products

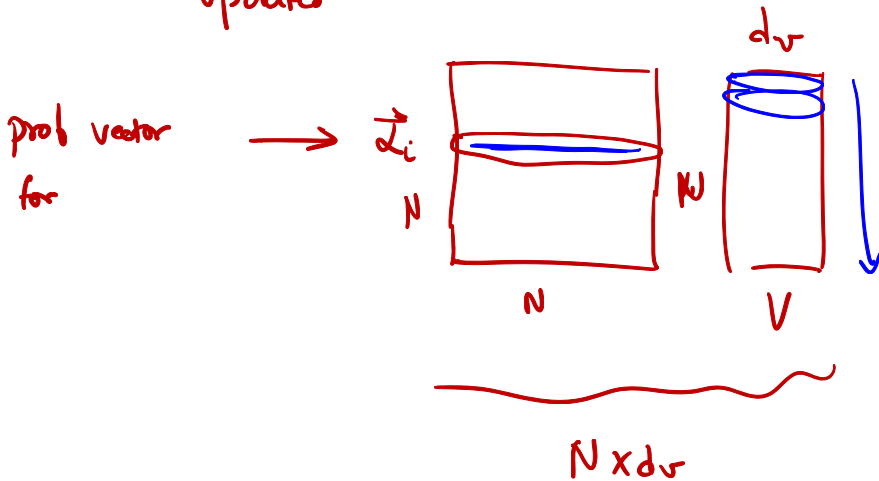


$A = \underset{\substack{\text{softmax} \\ \downarrow \\ \text{(rowwise)}}}{\left(\frac{Q \cdot K^T}{\sqrt{d_k}} \right)}$

d_{ij}
 $O(N^2)$

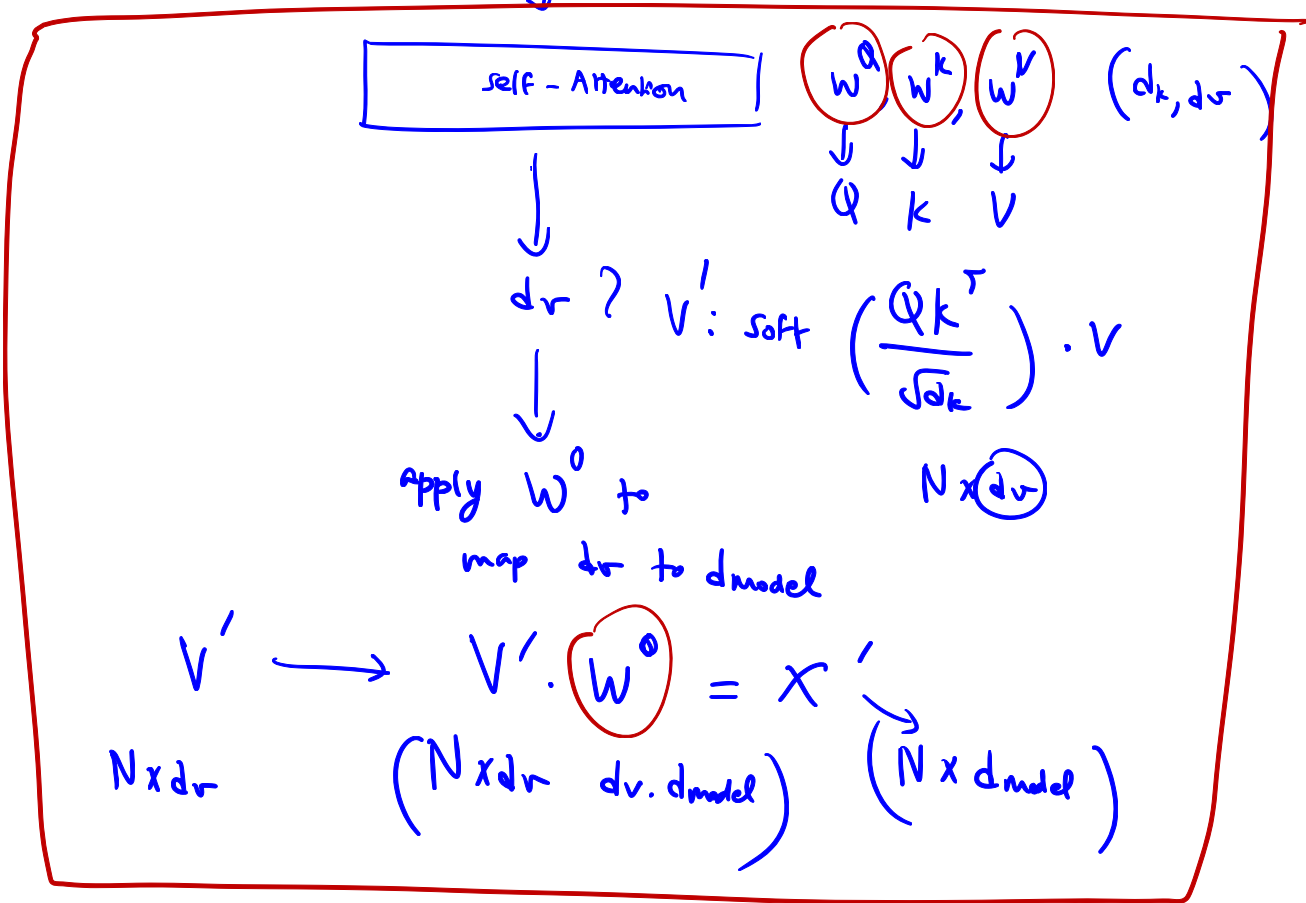
V : $A \cdot V$

updated

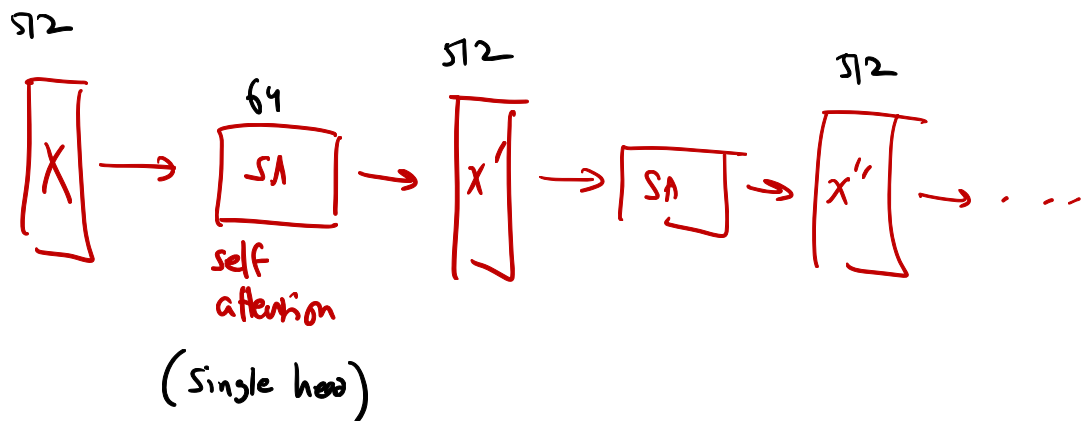


$$V' = \underbrace{\text{Softmax} \left(\frac{Q \cdot k^T}{\sqrt{d_k}} \right)}_A \cdot V$$

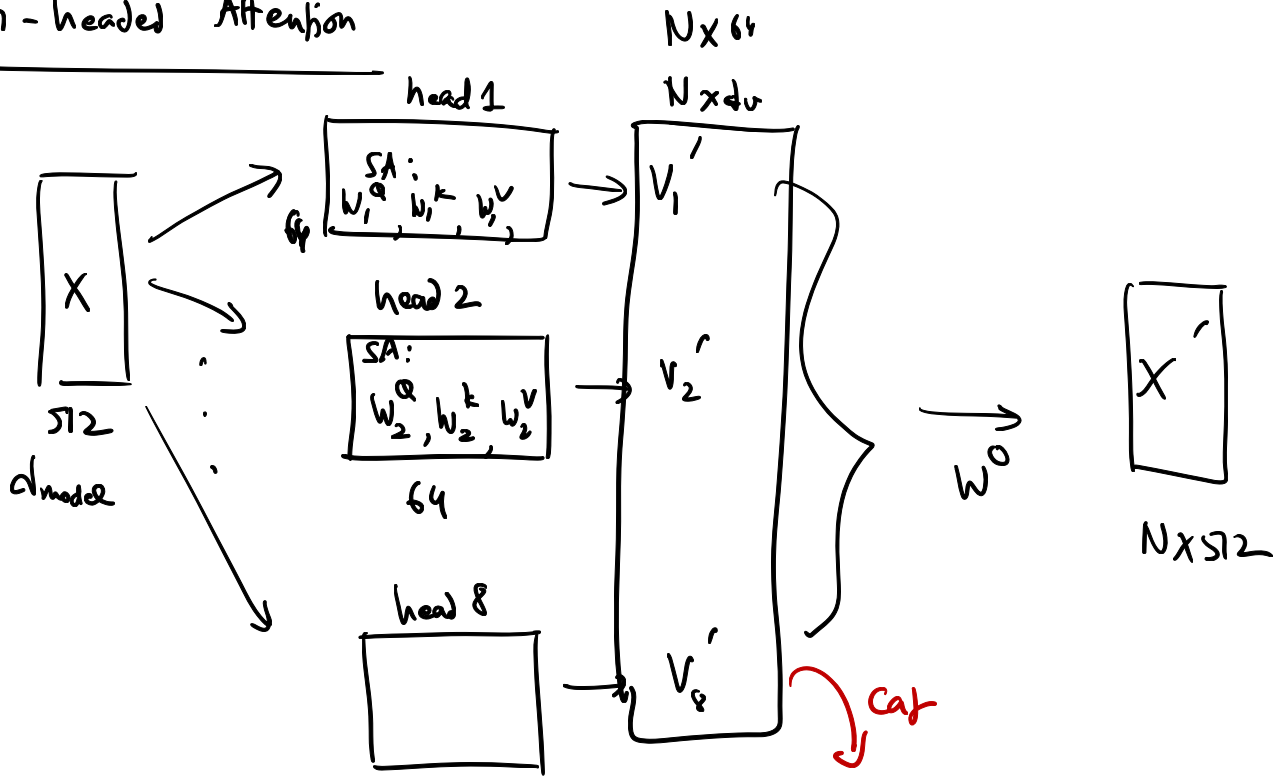
$$X \quad \boxed{[x_1 \quad x_2 \quad \dots \quad x_n]} \quad \underline{\underline{d_{model}}}$$



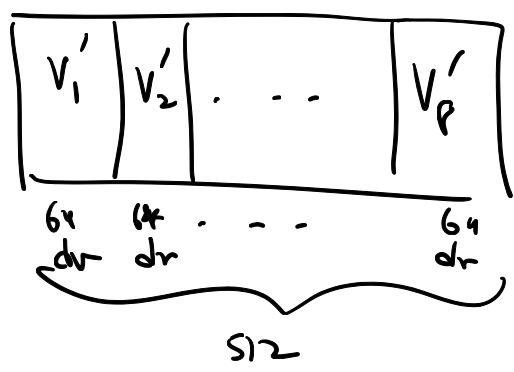
self-attention
 layer 0
 W
 $W^Q \quad W^K \quad W^V$



Multi-headed Attention



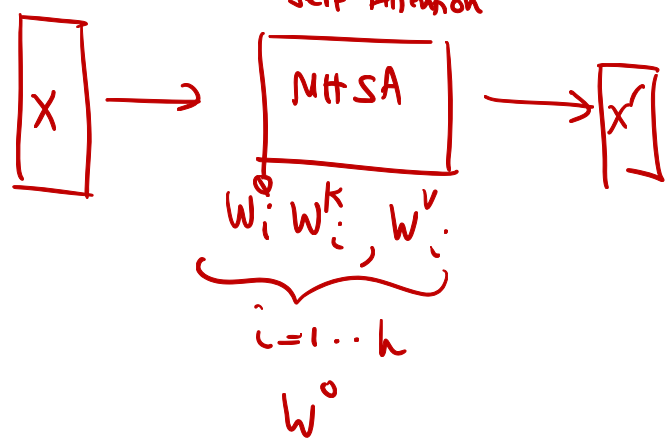
$C^V : N$



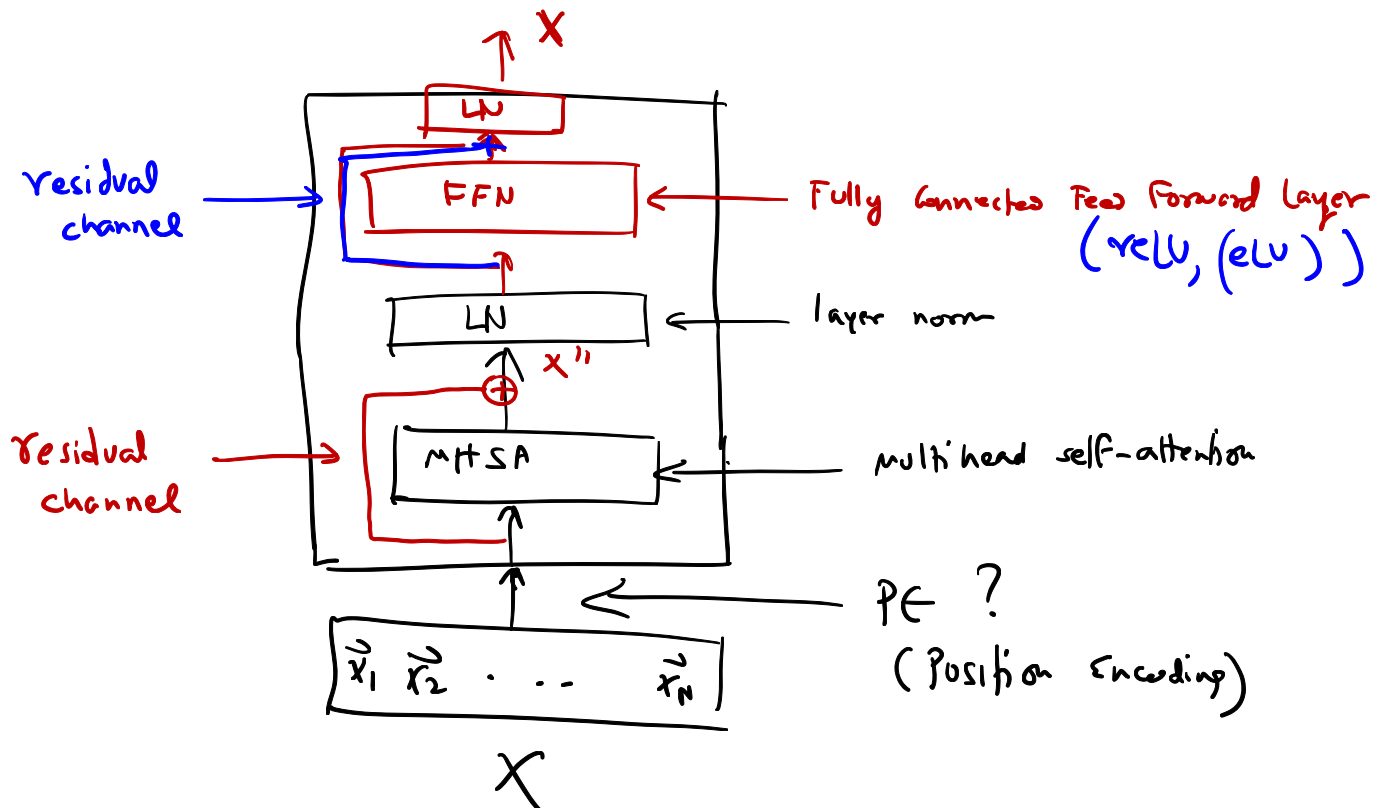
$$X' = C^V \cdot W^O$$

$(N \times h \cdot d_v) \quad (h \cdot d_v \times d_{model})$
 $\underline{\underline{8 \cdot 64 \times 512}}$

Multi-head self Attention

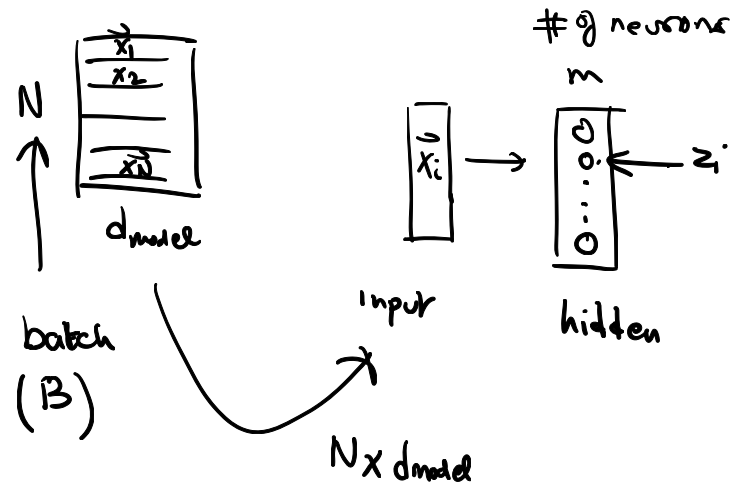


Transformer Block



Residual channel: $X'' = X + X' = \text{MHA}(X)$

Batch - Norm



Layer Norm

$$z_i = \alpha_i \left(\frac{z_i - \mu_1}{\sigma_1} \right) + \beta_i$$

$$\mu_1 = \frac{1}{m} \left(z_1(x_1) + z_2(x_1) + \dots + z_m(x_1) \right)$$

$$\sigma_1 = \sqrt{\frac{1}{m} \sum_i (z_i(x_1) - \mu_1)^2}$$

$$z_i = \frac{z_i - \mu_i}{\sigma_i} \quad (\text{z-score})$$

$$\rightarrow \mu_i = \frac{1}{N} \left(z_i(x_1) + z_i(x_2) + \dots + z_i(x_N) \right)$$

mean response across the batch

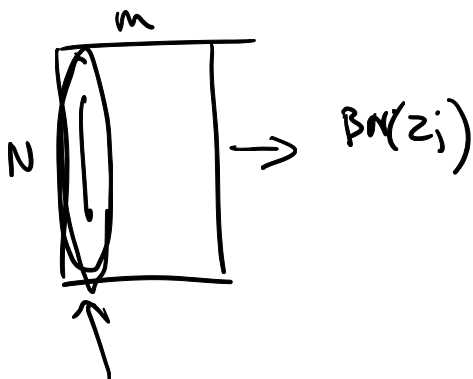
$$\sigma_i^2 = \frac{1}{N} \left((z_i(x_1) - \mu_i)^2 + \dots + (z_i(x_N) - \mu_i)^2 \right)$$

$$\sigma_i = \sqrt{\sigma_i^2}$$

$$z_i = \alpha_i \left(\frac{z_i - \mu_i}{\sigma_i} \right) + \beta_i$$

↑ gain ↑ bias

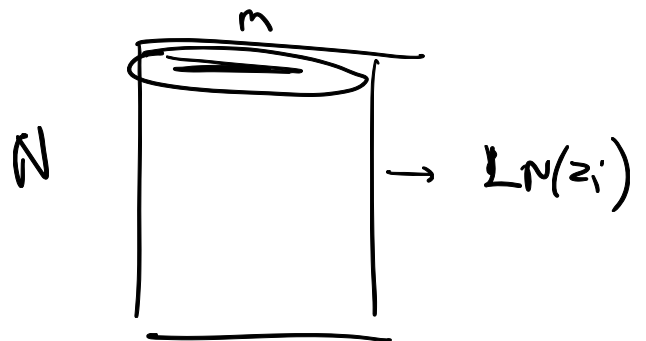
Batch Norm



avg across the batch elements

$$\mu = \frac{1}{N} \sum ()$$

layer norm



avg across the layer (m)

$$\frac{1}{m} \sum ()$$

Word2vec / Pnt2vec

Vocab + Counts

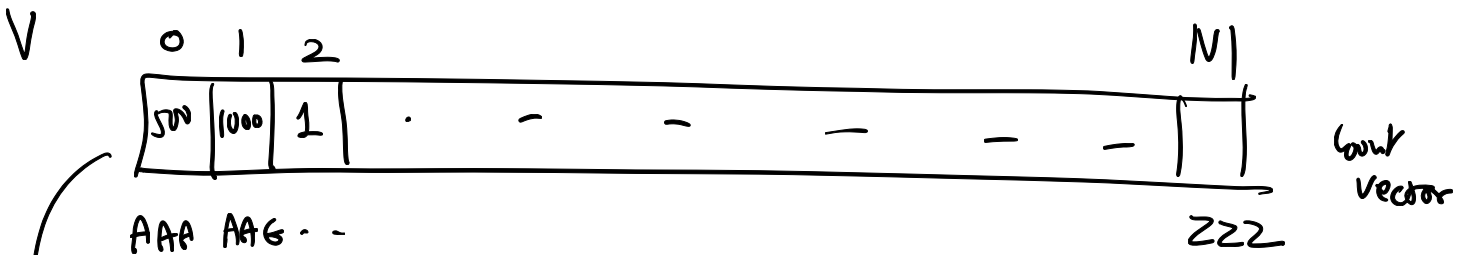
k-mer $k=3$



do k-way sliding window

Collect the vocab + count in (defaultdict)

Negative sampling



$|V|=4$

