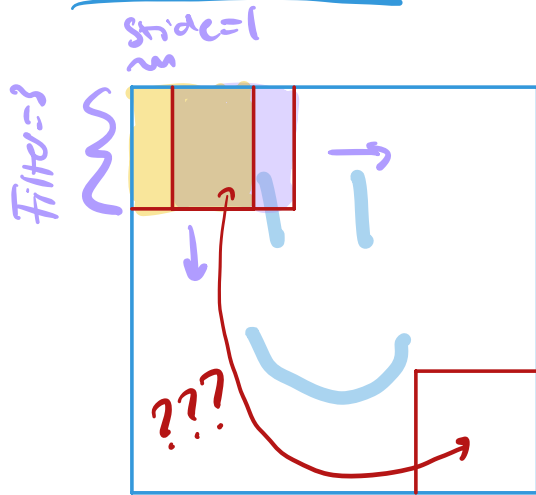


Convolution

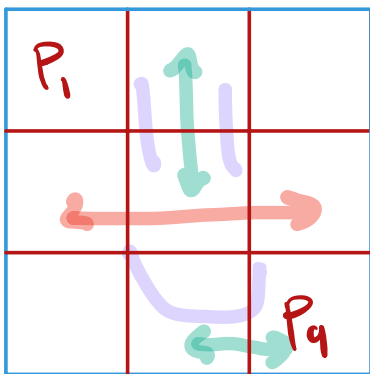


Sliding filter approach

- Local computation focused on a small piece of image.
- No learned relation btwn different located filters

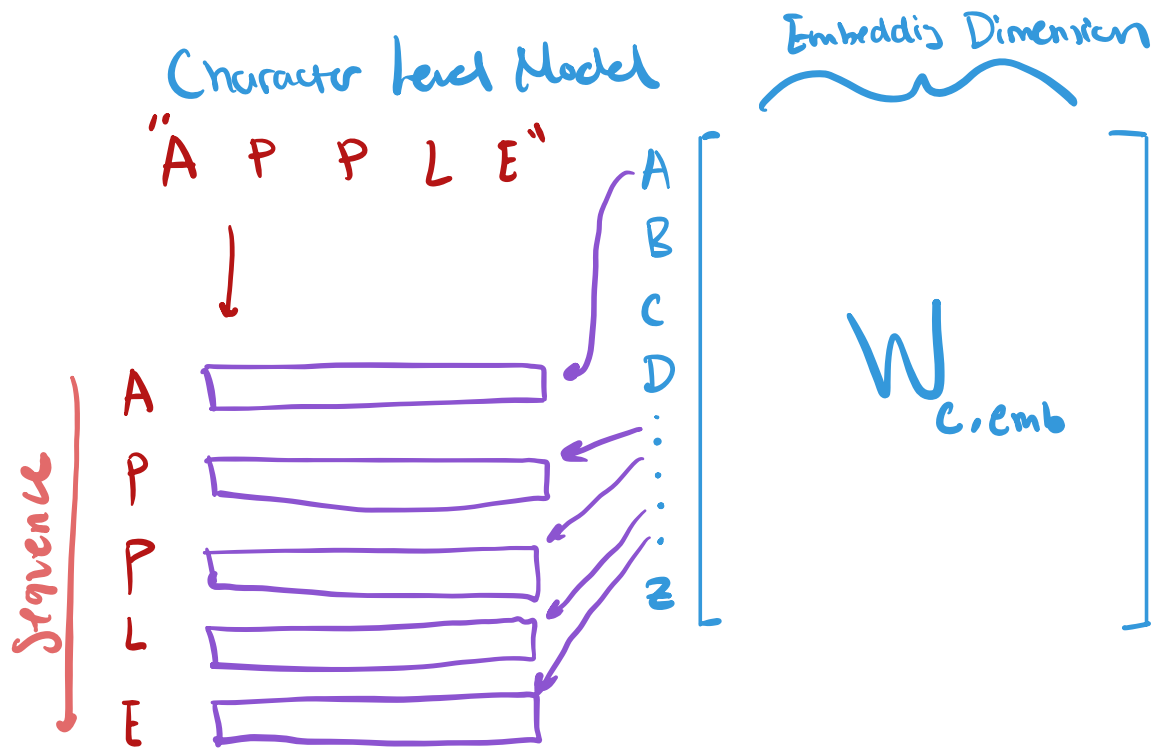
ViT

- capture global image relations

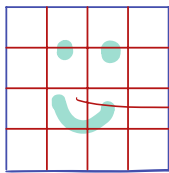


Q. How is patch 1 related to patch 9?

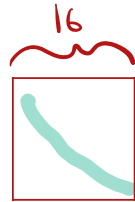
The typical NLP Tasks (Tokenizing/Embedding)



Vision Transformer → How do we convert an image to an embedded sequence.

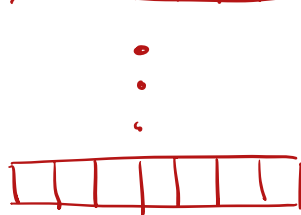
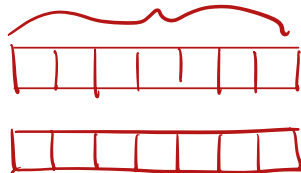


The only "extra" thing the ViT paper did. Afterwards, almost the same as Orig transformer paper.



16

768

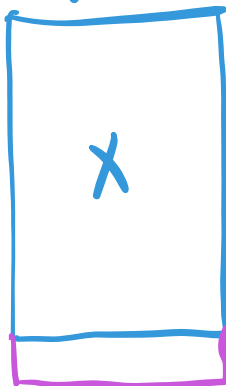


patches + 1
(seq-len)

Concatenation
CLS Token



768



P+1

+

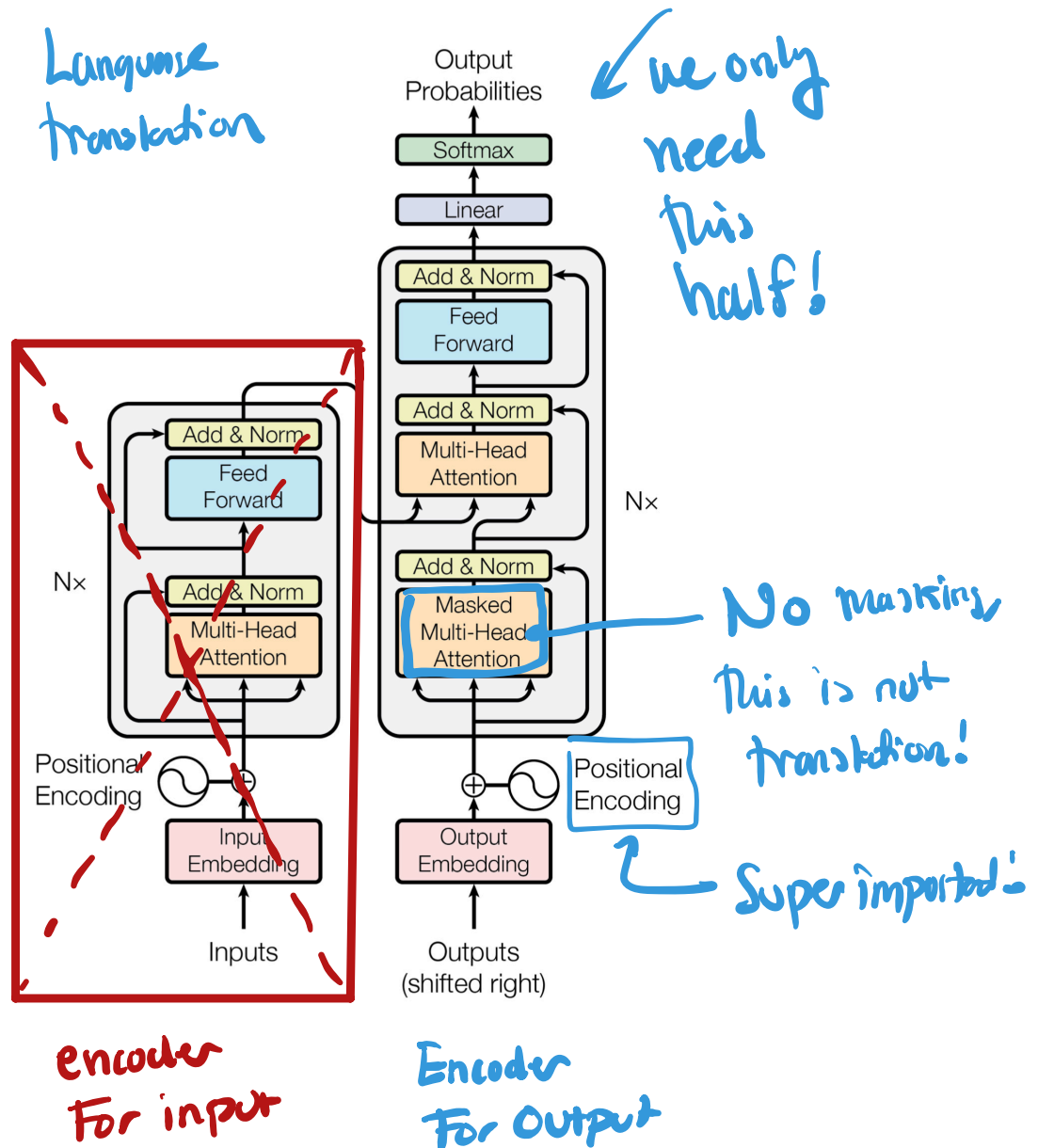
768

Position
embedding

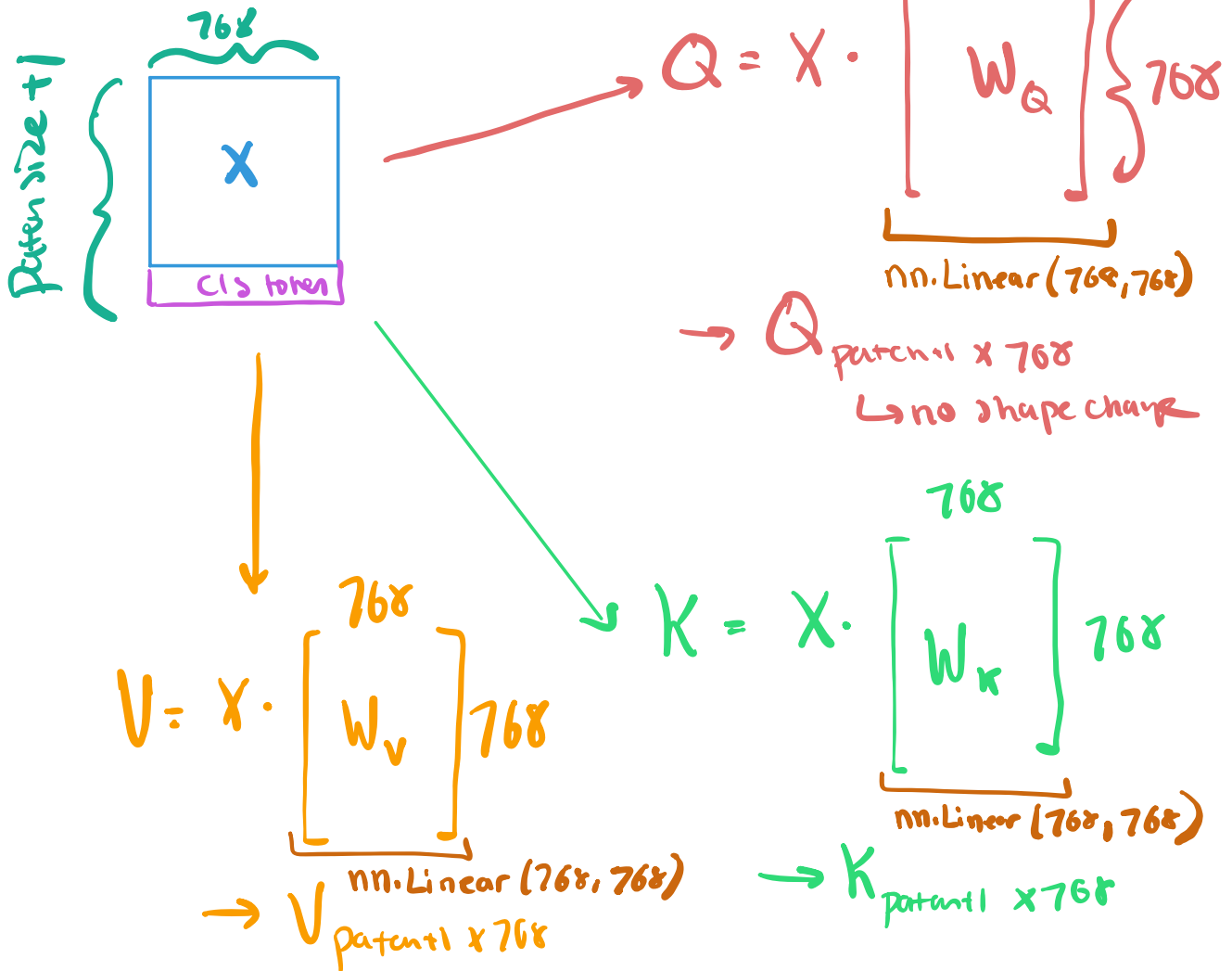
P+1

↑ inject positional
info, otherwise
permutation
invariant

What exactly is a transformer?

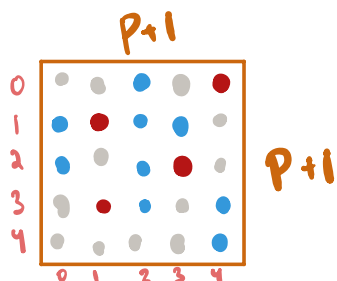


Single Attention Head



Attention Matrix (Fancy weighted average)

$$Q_{p+1, 768} \cdot K_{768, p+1}^T \rightarrow \text{Attention mat.}$$



$A_{p+1, p+1}$

$$\frac{A_{p+1}}{\sqrt{768}}$$

← divide by sqrt head size
→ scaled attention.

Variance of A will scale w/
head size so we want to
normalize

$$\begin{aligned} \text{if } \text{Var}(A) &= 768 \cdot \text{Var}(A^*) \\ \rightarrow \text{Var}\left(\frac{A}{\sqrt{768}}\right) &= \left(\frac{1}{\sqrt{768}}\right)^2 \cdot 768 \cdot \text{Var}(A^*) \\ &= \text{Var}(A^*) \\ &\approx \text{Var}(Q), \text{Var}(K) \end{aligned}$$

$\text{Softmax}(A_{p+1,p+1}^*) \rightarrow$ scale to probability vector
between $[0,1]$

$$A_{p+1,p+1}^{*s} \cdot V_{p+1 \times 768} \rightarrow O_{p+1, 768}$$

same shape!

$$A(Q, K, V) = \text{Softmax}\left(\frac{Q \cdot K^T}{\sqrt{d_k}}\right) V$$

$$P+1 \begin{bmatrix} \overset{P+1}{A^{*S}_{P+1, P+1}} \end{bmatrix} \cdot \begin{bmatrix} \overset{768}{V_{P+1, 768}} \end{bmatrix} \quad P+1$$

Values are a projection of the original data X

That encode each patch w/ an embedding vector.

We want to do a weighted average of those

embeddings.

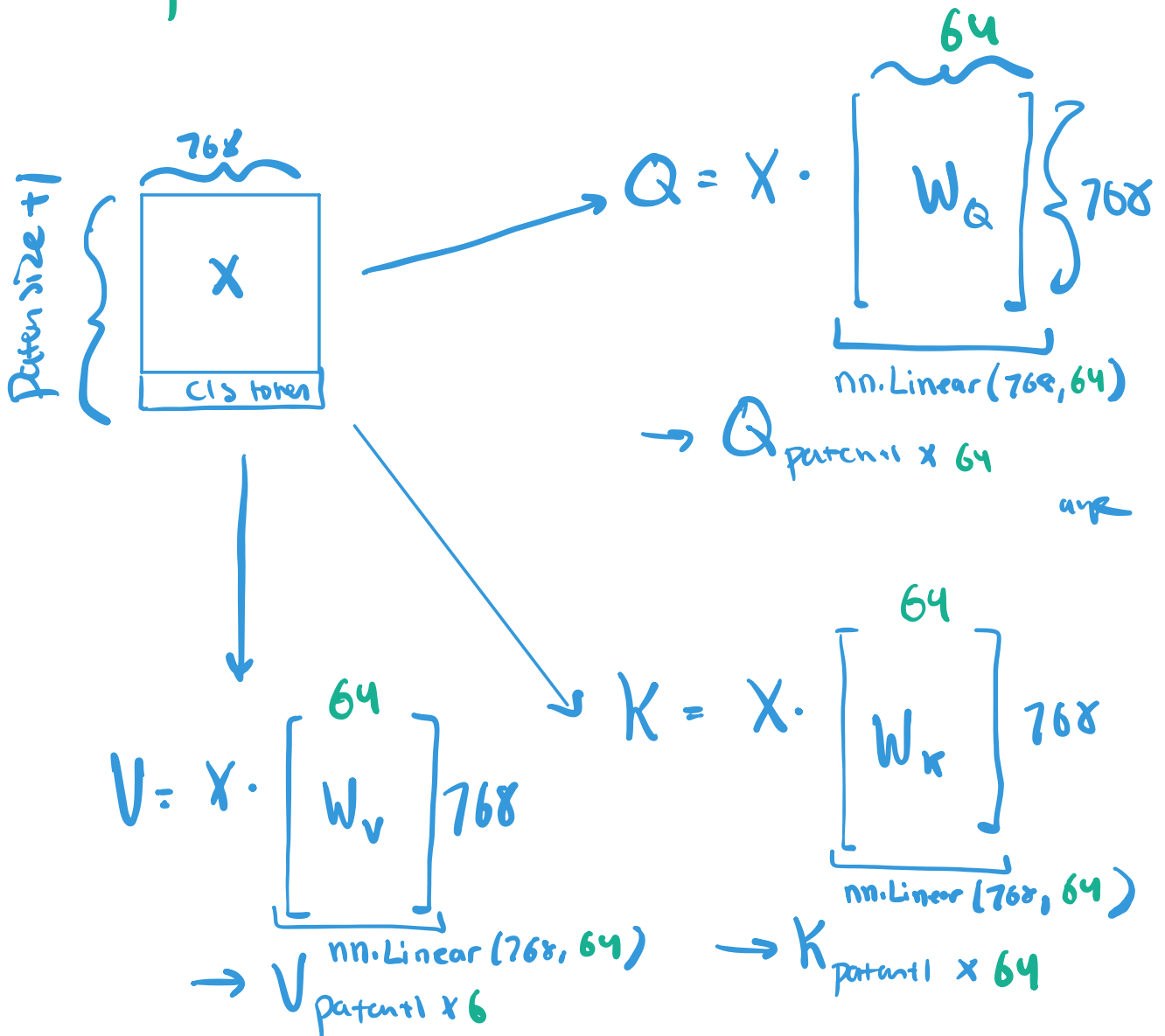
$$P+1 \begin{bmatrix} \overset{P+1 (3)}{0.2 \quad 0.3 \quad 0.5} \\ \vdots \\ \vdots \end{bmatrix} \begin{bmatrix} \overset{768}{-e_1-} \\ -e_2- \\ -e_3- \end{bmatrix} \quad P+1 (3)$$

generated by W_v proj

$$\rightarrow \begin{bmatrix} 0.2e_1 + 0.3e_2 + 0.5e_3 \\ \vdots \\ \vdots \end{bmatrix} \leftarrow \text{weighted average.}$$

Multiheaded Attention

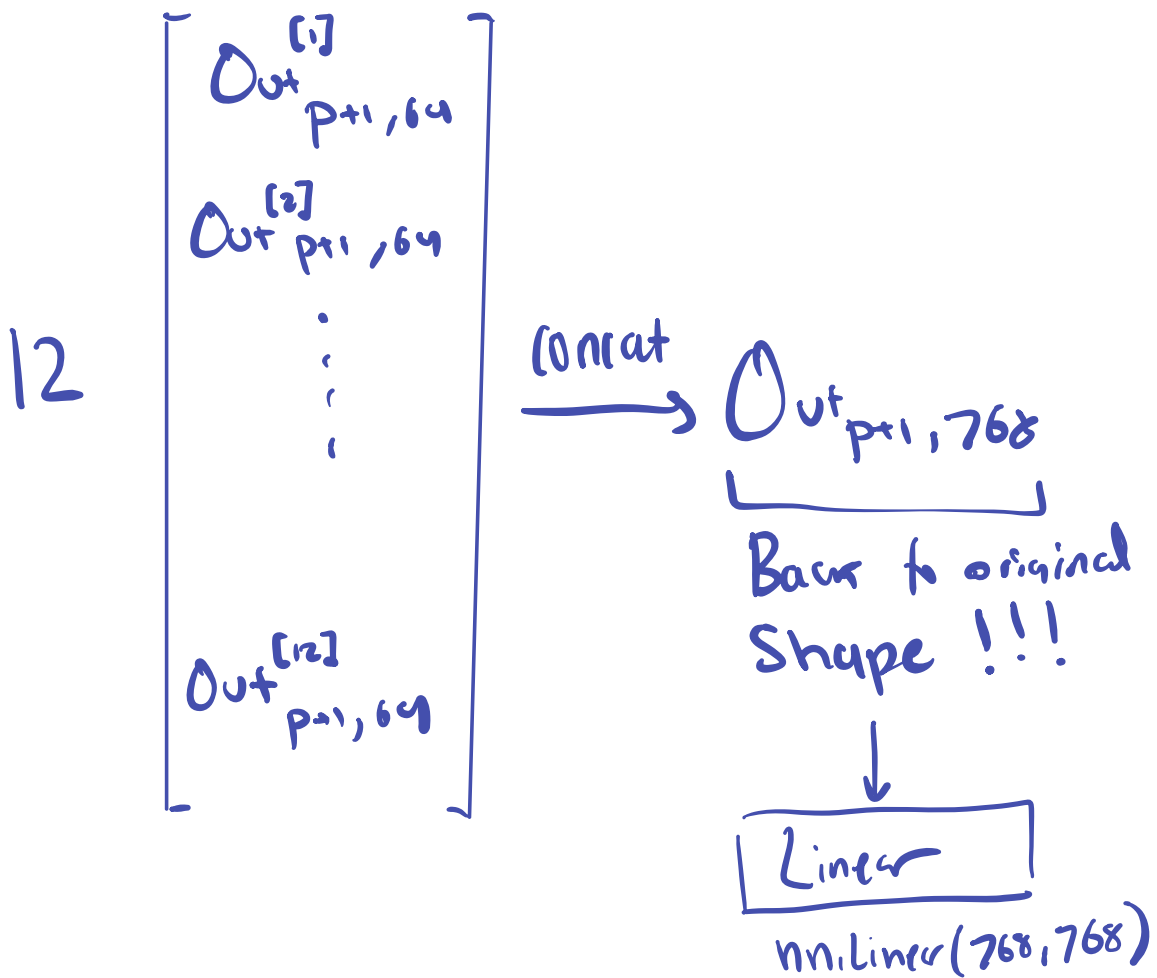
if we start with embedding = 768 and we want 12 heads $\rightarrow 768/12 = 64$ dim per head.



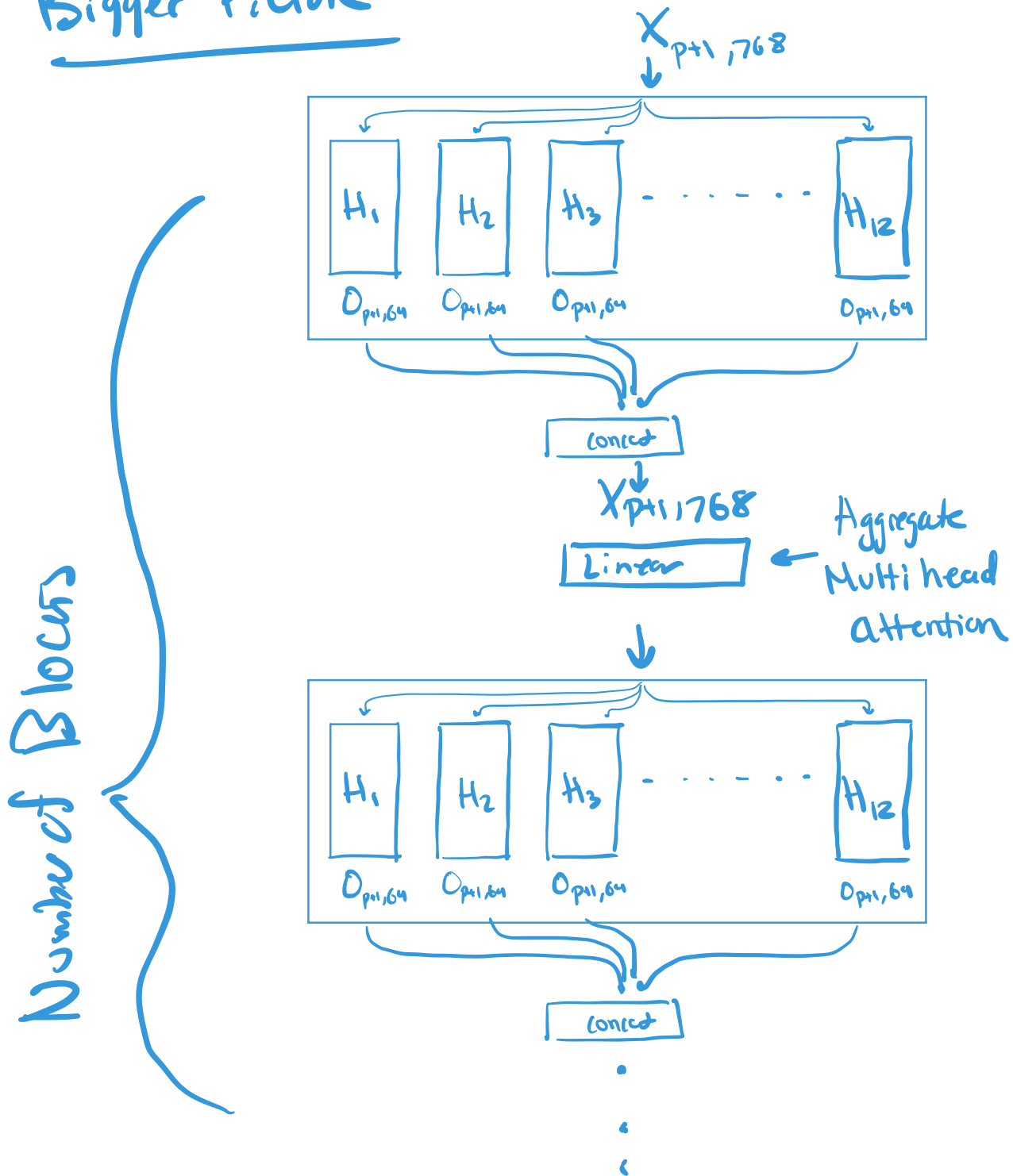
All previous calc. For single head identical
so a single head returns:

$$Out_{patten+1, \underline{\underline{64}}}$$

But we have 12 heads so output will
really be:



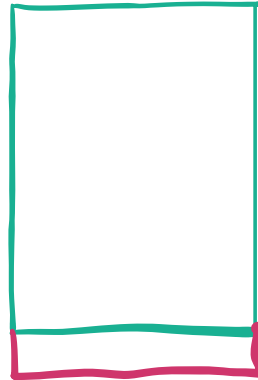
Bigger Picture



$$\text{Total Attention} = \text{num blocks} \times \text{num heads}$$

Classifier

$O_{p+1, 768}$ →



Want CLS token to accumulate info from other patterns

← CLS token
1x768

→ nn.Linear(768, n-classes)