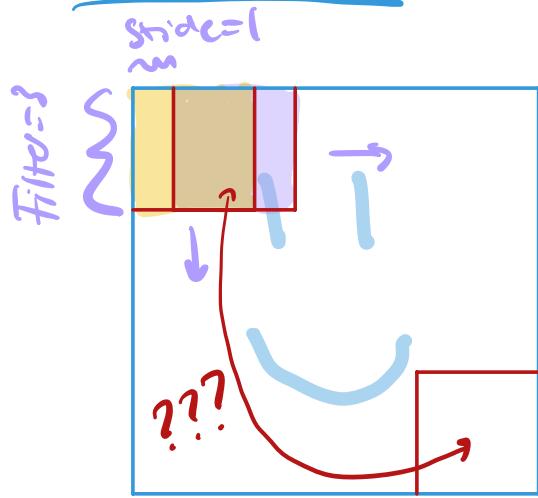


Convolution

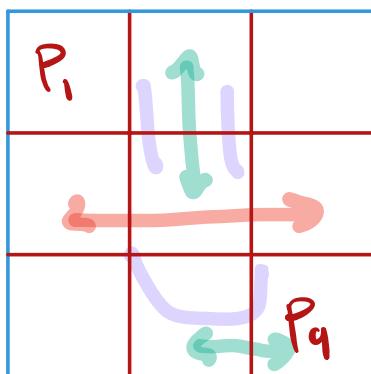


Sliding filter approach

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

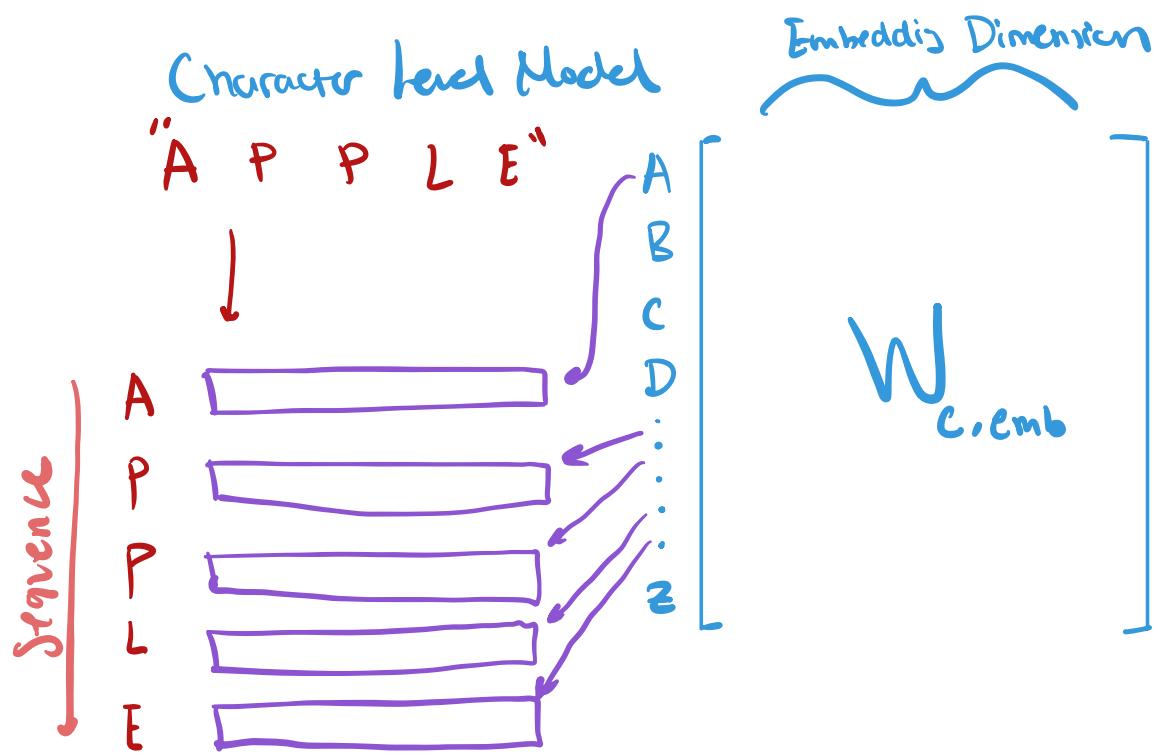
ViT

→ capture global image relations

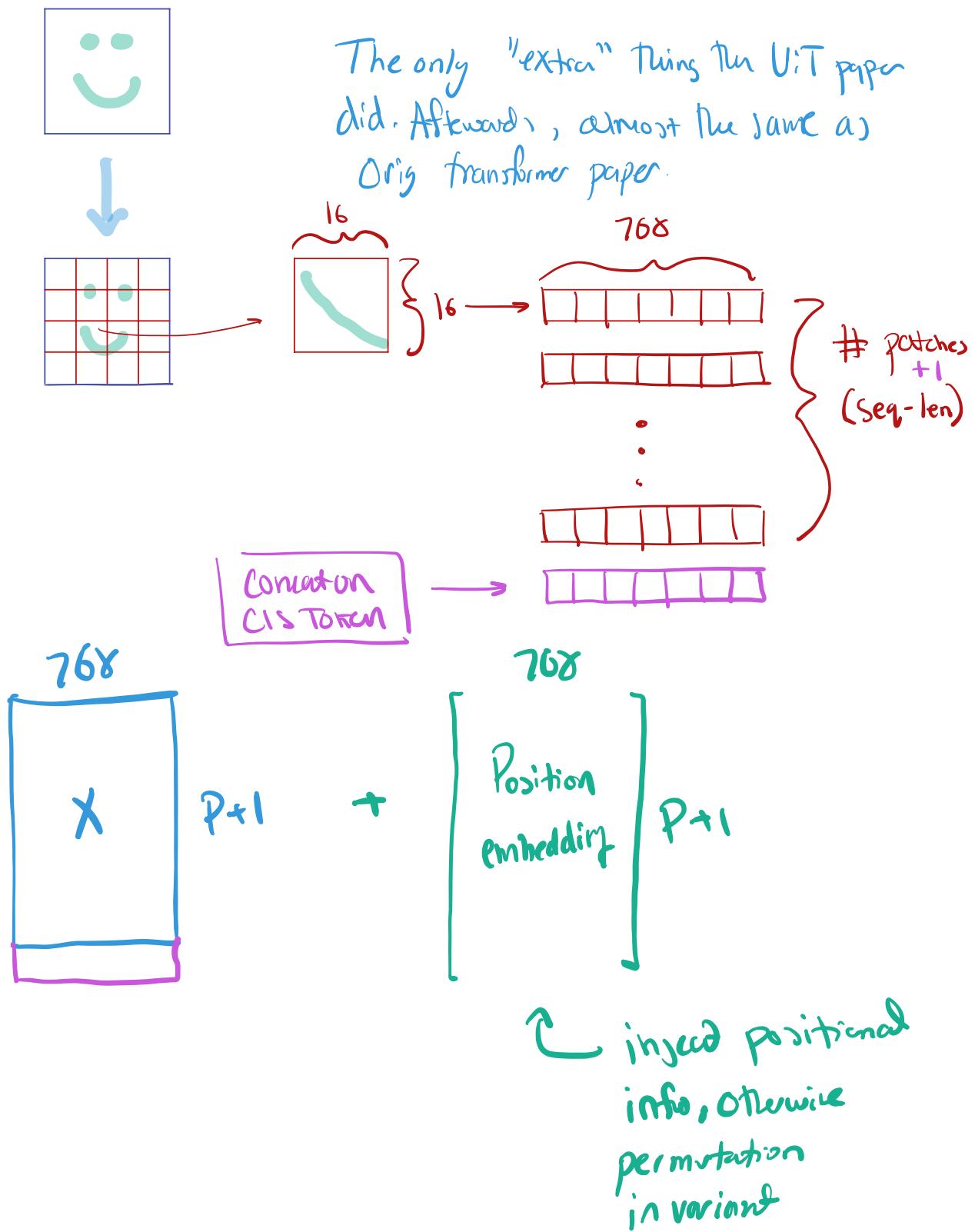


Q. How is patch 1 related to patch q?

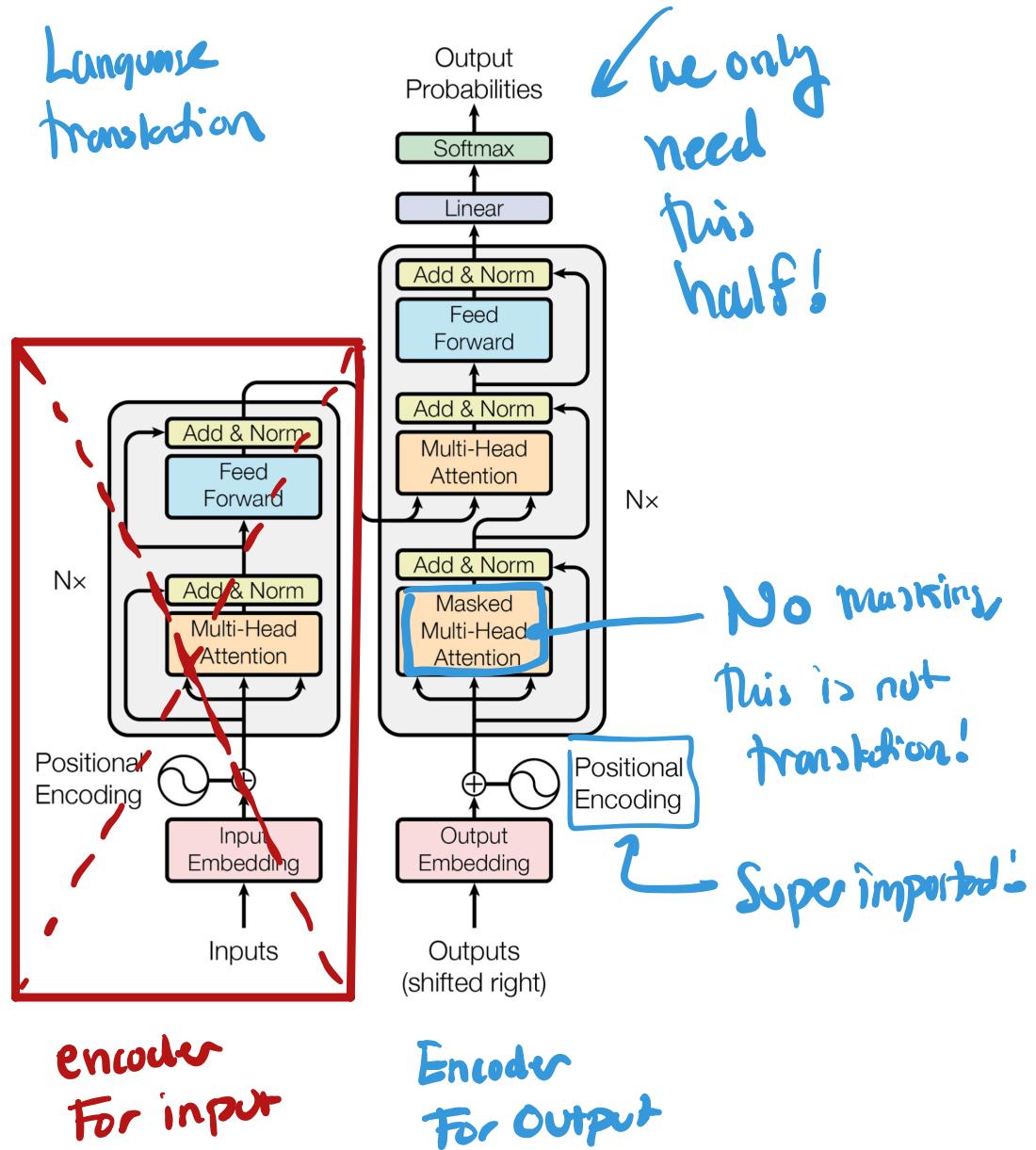
The typical NLP Task (Tokenizing/Embedding)

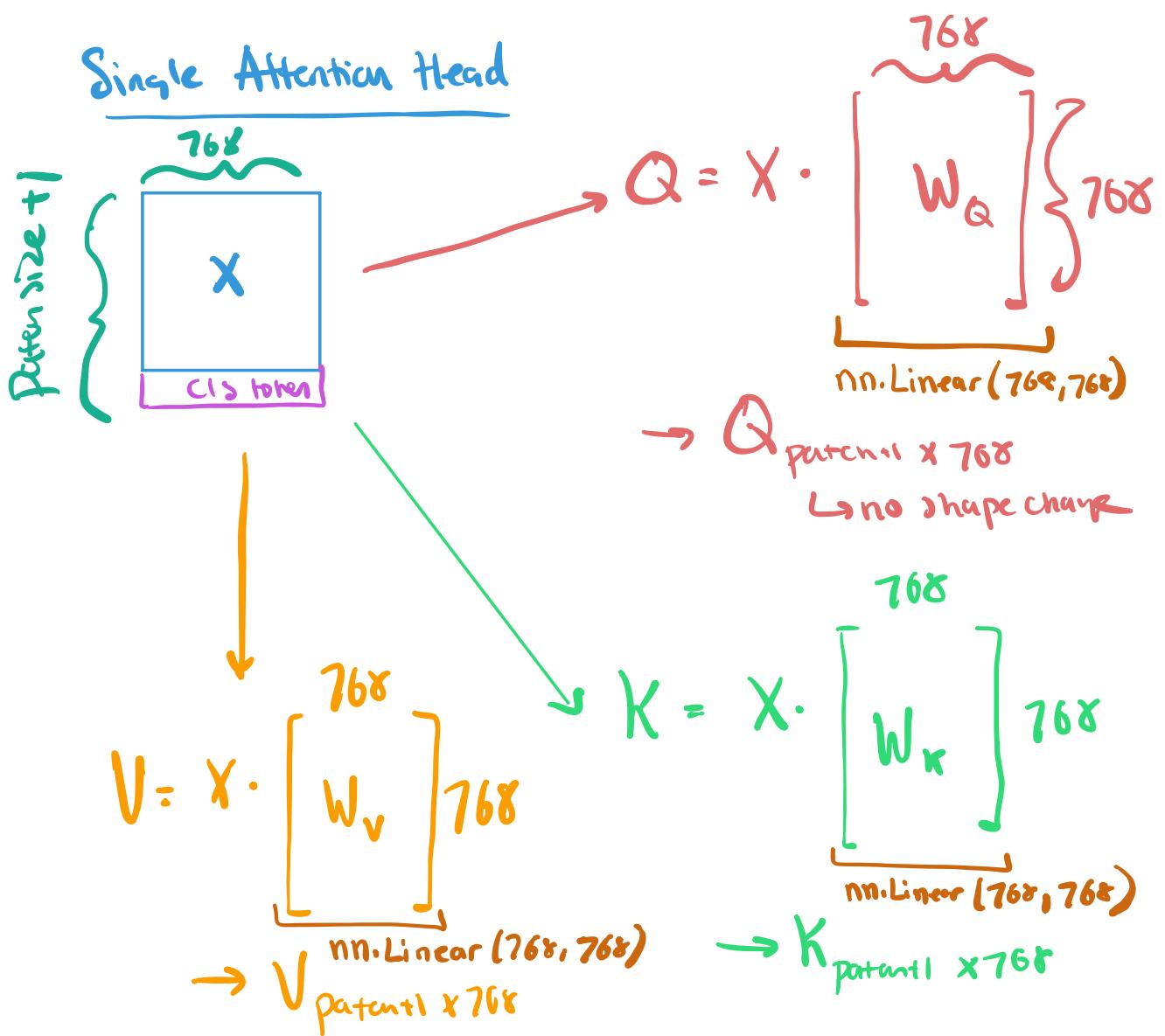


Vision Transformer → How do we convert a image to an embedded sequence.



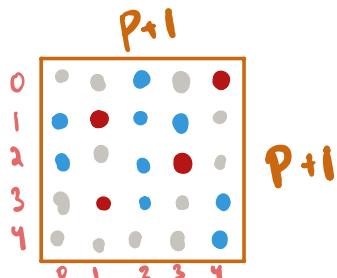
What exactly is a transformer?





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

if $\text{Var}(A) = 768 \cdot \text{Var}(A^*)$

$$\rightarrow \text{Var}\left(\frac{A}{\sqrt{768}}\right) = \left(\frac{1}{\sqrt{768}}\right)^2 \cdot 768 \cdot \text{Var}(A^*)$$

$$= \text{Var}(A^*)$$

$$\approx \text{Var}(Q), \text{Var}(K)$$

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

$$A_{p+1, p+1}^{*s} \cdot V_{p+1 \times 768} \rightarrow \underbrace{\text{Out}_{p+1, 768}}_{\text{same shape!}}$$

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

$$P+1 \begin{bmatrix} P+1 \\ A_{P+1, 768}^* \end{bmatrix} \cdot \begin{bmatrix} \vdots \\ V_{P+1, 768} \end{bmatrix} \quad P+1$$

Values are a projection of the original data X

that encode each patch w/ a embeddings vector.
We want to do a weighted average of those

embeddings.

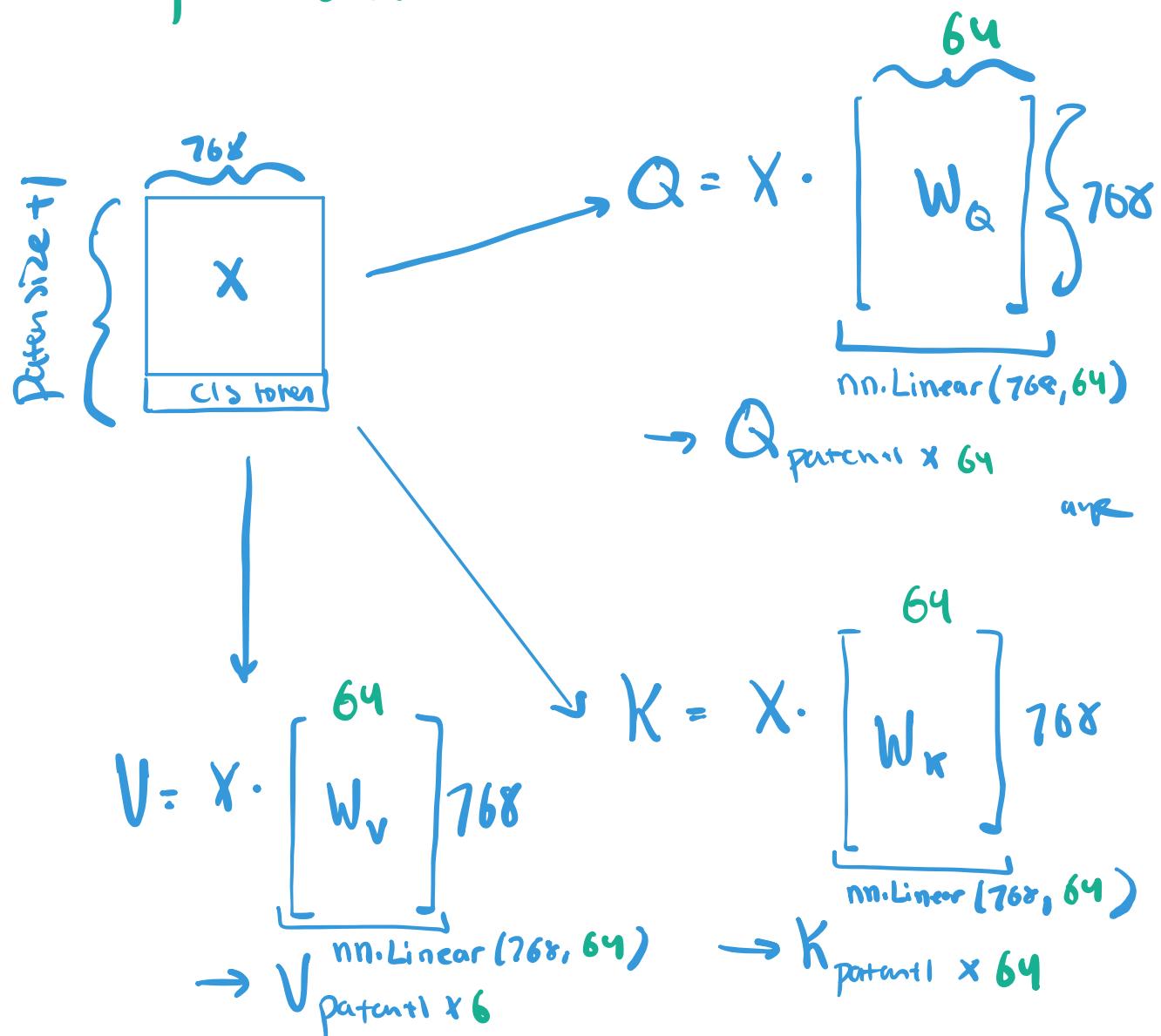
$$P+1 \begin{bmatrix} P+1 & (3) \\ \begin{bmatrix} 0.2 & 0.3 & 0.5 \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \end{bmatrix} \end{bmatrix} \begin{bmatrix} 768 \\ -e_1- \\ -e_2- \\ -e_3- \end{bmatrix} \quad P+1 \begin{bmatrix} (3) \end{bmatrix}$$

generated by W_v proj

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

Multiheaded Attention

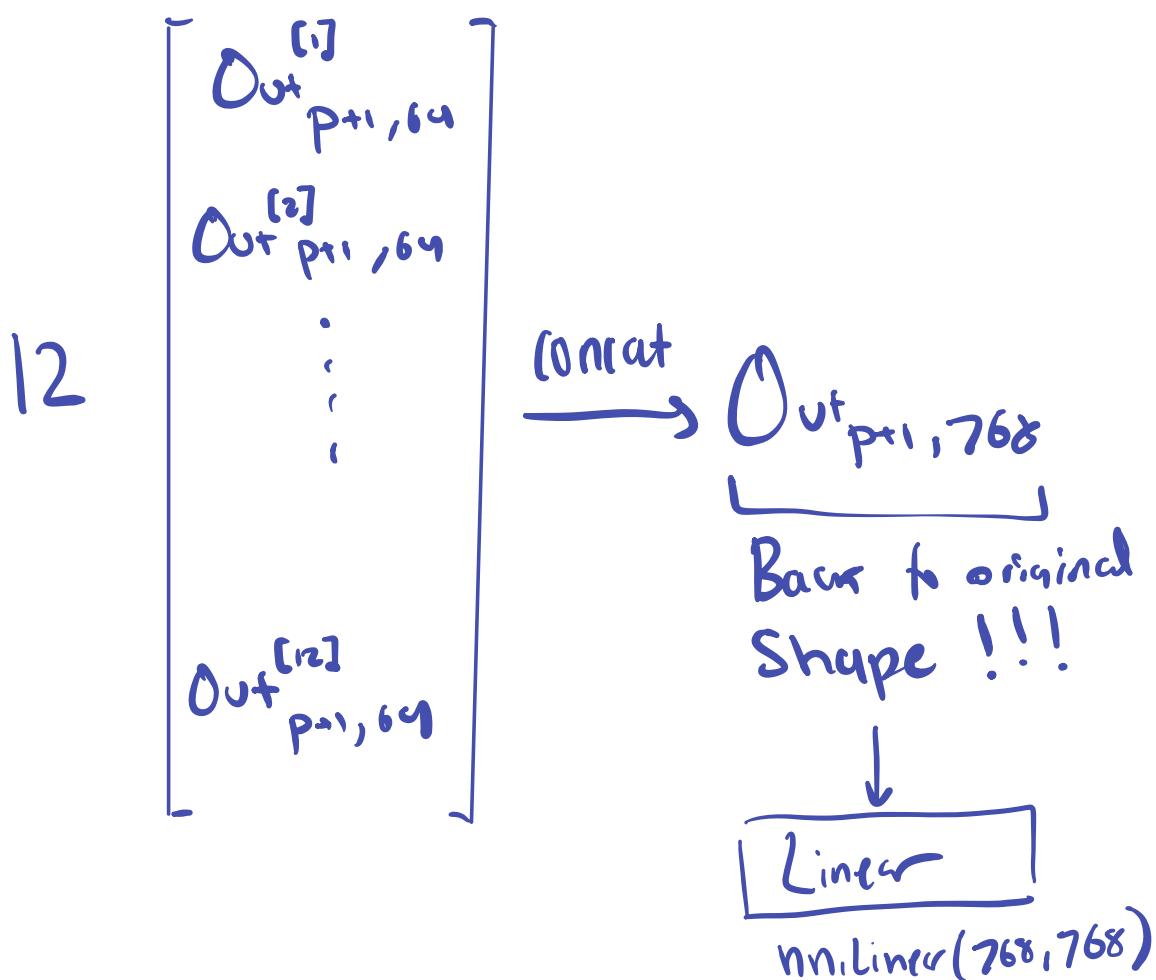
if we start with embedding = 768 and we want 12 heads $\rightarrow 768/12 = 64$ dim per head.



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

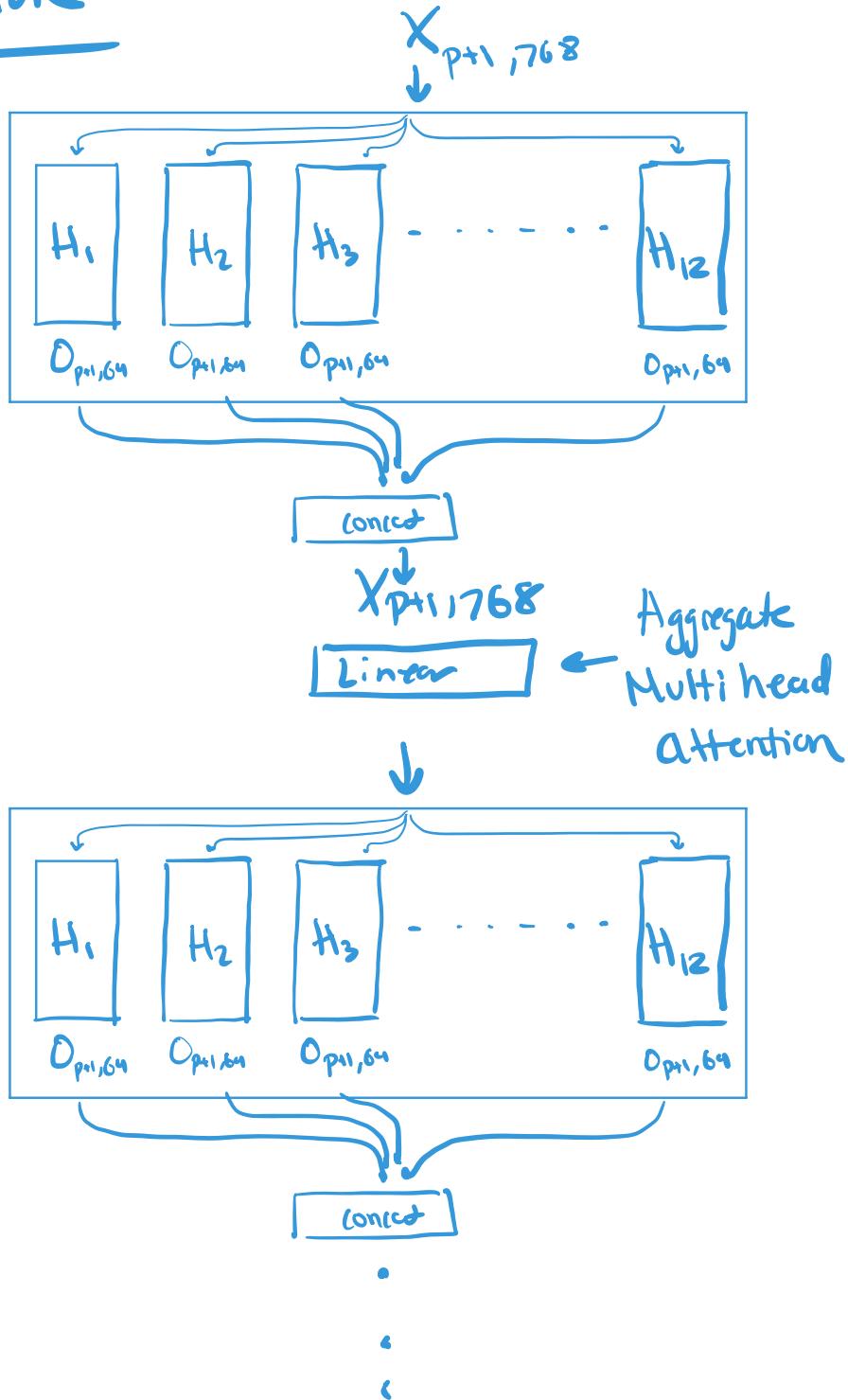
Out_{p+1, 64}

But we have 12 heads so output will
really be:



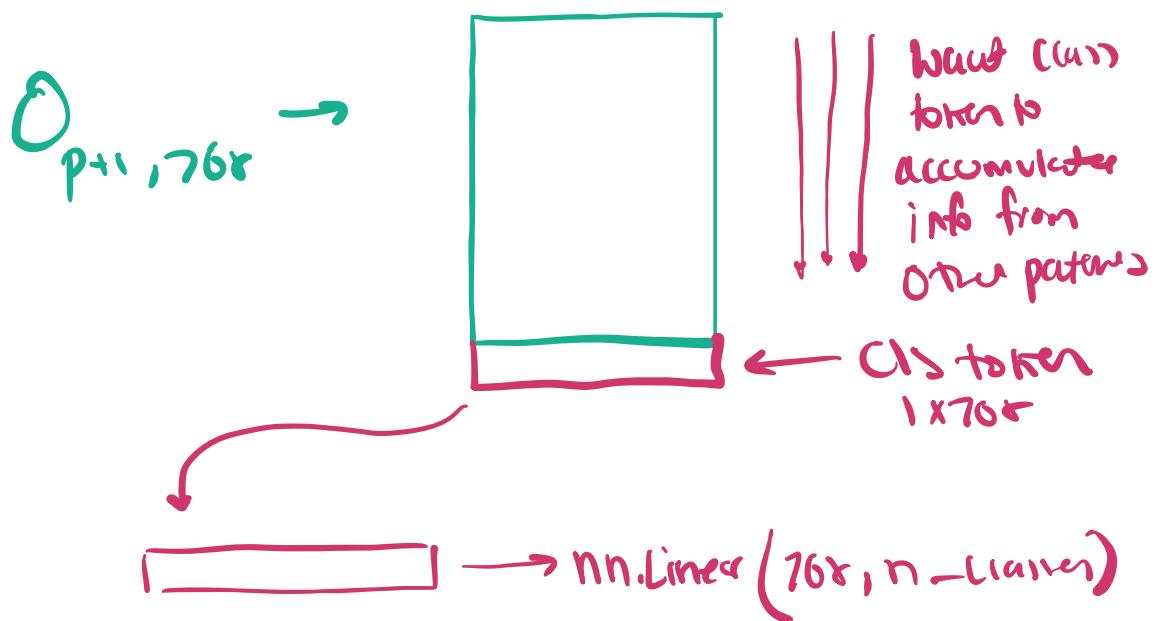
Bigger Picture

Number of Blocks

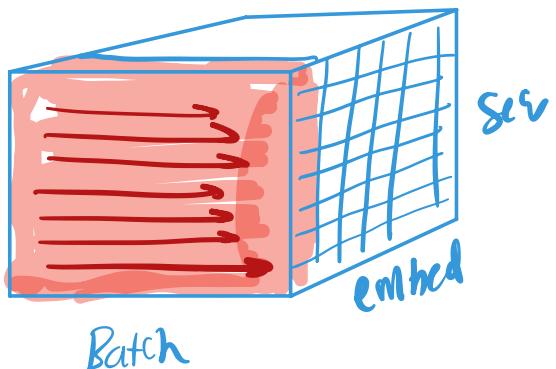


$$\text{Total Attention} = \text{num blocks} \times \text{Num heads}$$

Classifier

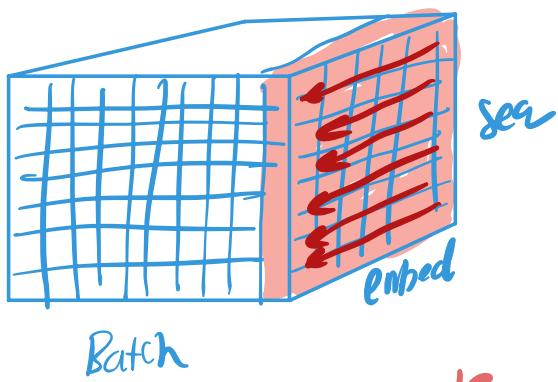


Batch Norm



Batch
calc μ and σ^2
for every pixel
across batch

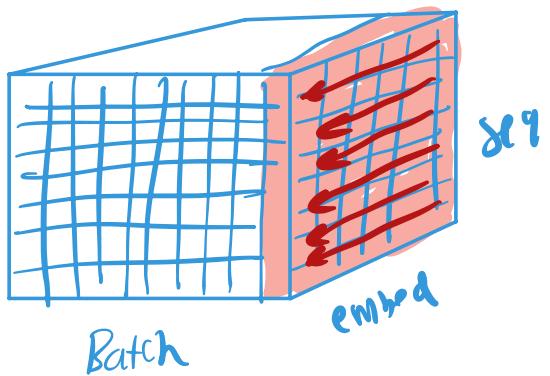
Layernorm



normalize each sample
individually
→ don't need to
store μ and σ
anymore!

Problems w/ batchnorm

- hard to use w/ different sequence lengths
- fails w/ small batch sizes
- hard to parallelize



Input to LN is
(batch, seq, embed)

Goal: Calculate μ, σ for
each embedding vector for
each sample in the batch

`nn.LayerNorm(embed_dim)`

↑
each embedding vector is
of length embed_dim, so
we have to let layer norm
what to expect

From docs: "If single integer is used, this module
will normalize over the last dimension
expected to be of this specific
size"