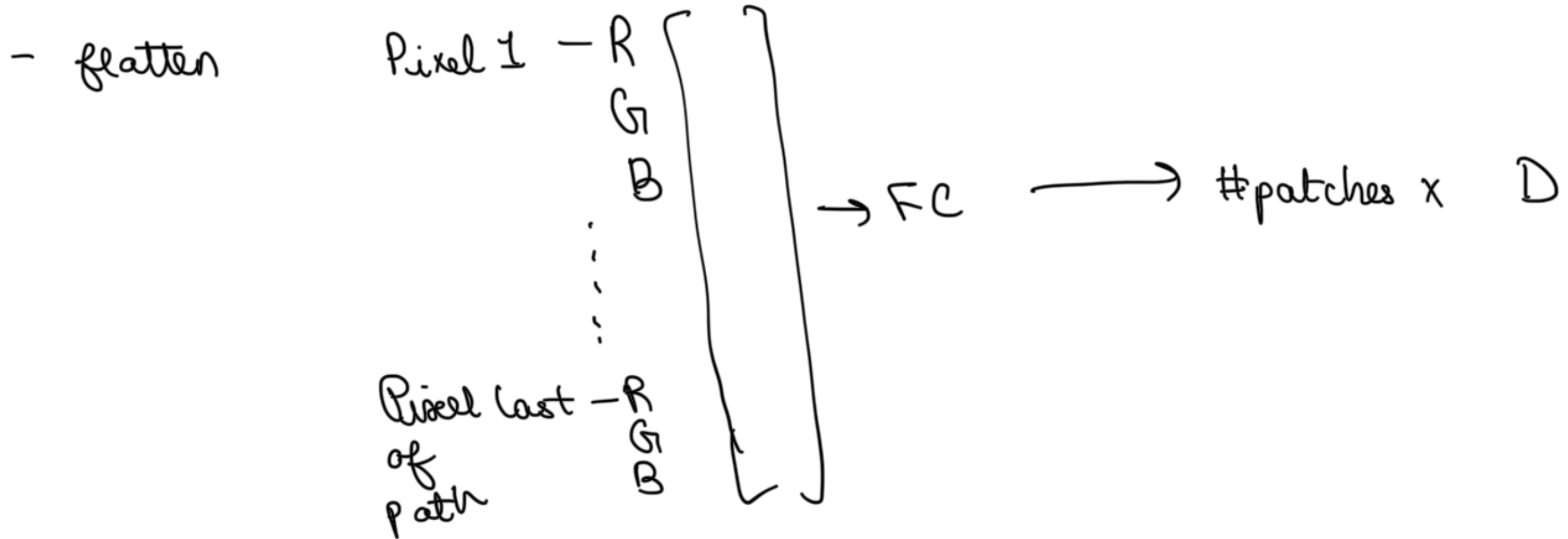
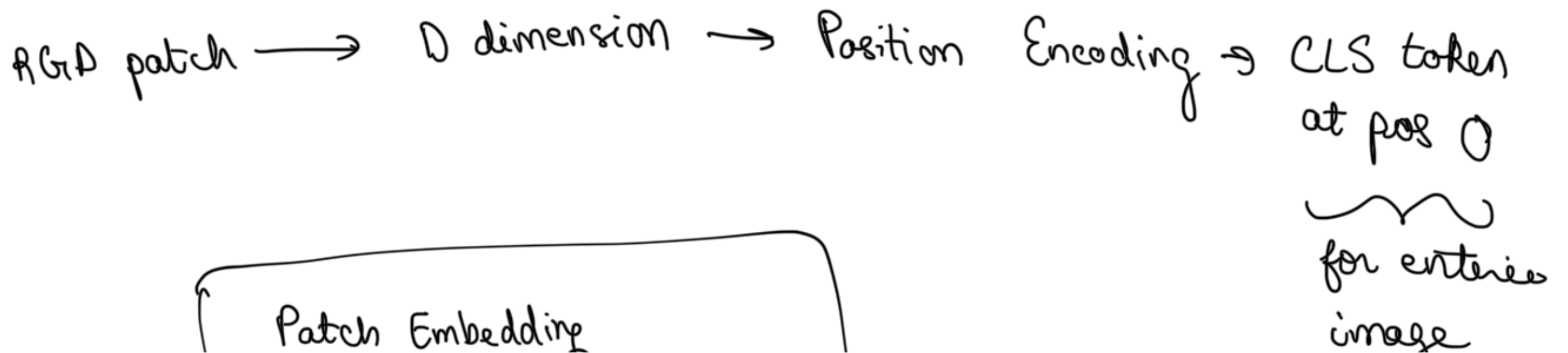


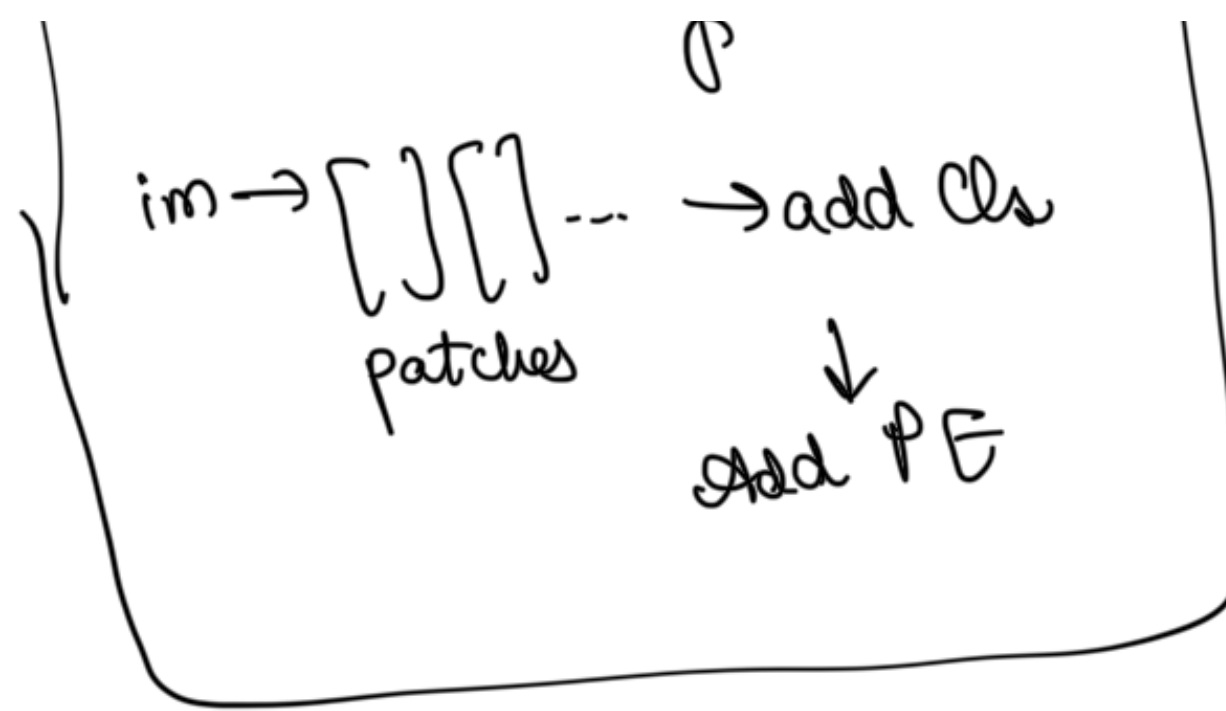


-  $\frac{H \cdot W}{p \cdot sq \times p \cdot sq} \times 3$  channels



#patches x channels  
x patch ht x patch width





Enops

eg.  $\frac{224 \cdot 224}{16 \cdot 16} = 196 \text{ patches}$

-16x16 words

Requirement : Find relevance b/w Patch1 and Patch 4

Existing representation might not cater to what's relevant

$w_s$

$w_k$

$w_v$

learn  
Criteria  
of Relevance

how to  
so that  
criteria

transform patch  
based on  
similar patches

$D \times \text{head dimension}$

give similar products

$$(1 \times D) \times (D \times \text{head dim})$$

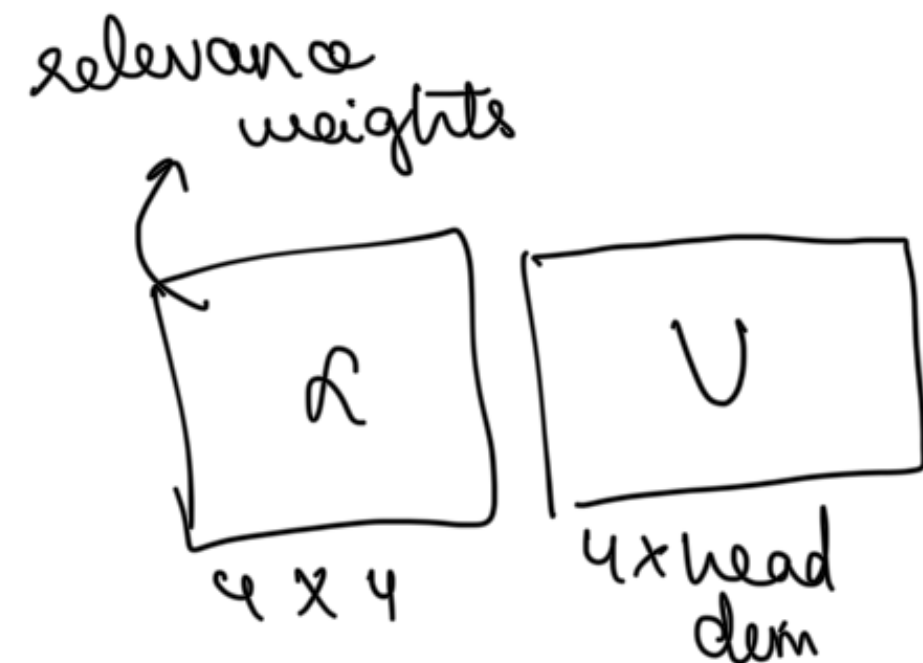
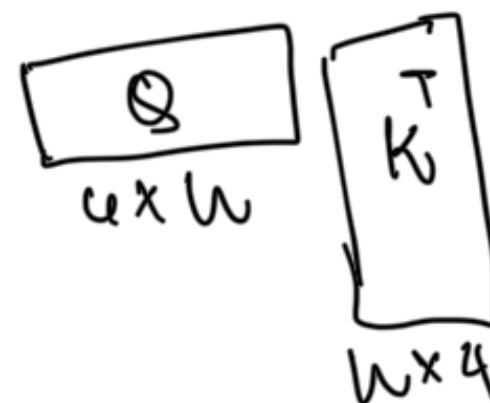
↑ Patches  
will be  
transformed to  
this dimension

$$d_i = \frac{Q K_i^T}{\sqrt{D_h}}$$

for patch 1, use  $Q$  to find which  
patches are most relevant to it

$$(d_1, d_2, \dots, d_{\# \text{patches}}) \leftarrow \text{softmax}$$

let #patches be 4



✓  
 $4 \times D$

$W_K$

"

$W_V$

$V$

$D \times \text{head dimension}$

$4 \times \text{head dim}$

patch

heads  
Query

$d_{i,j}$   
Query  $i$  Key  $j$

$=$

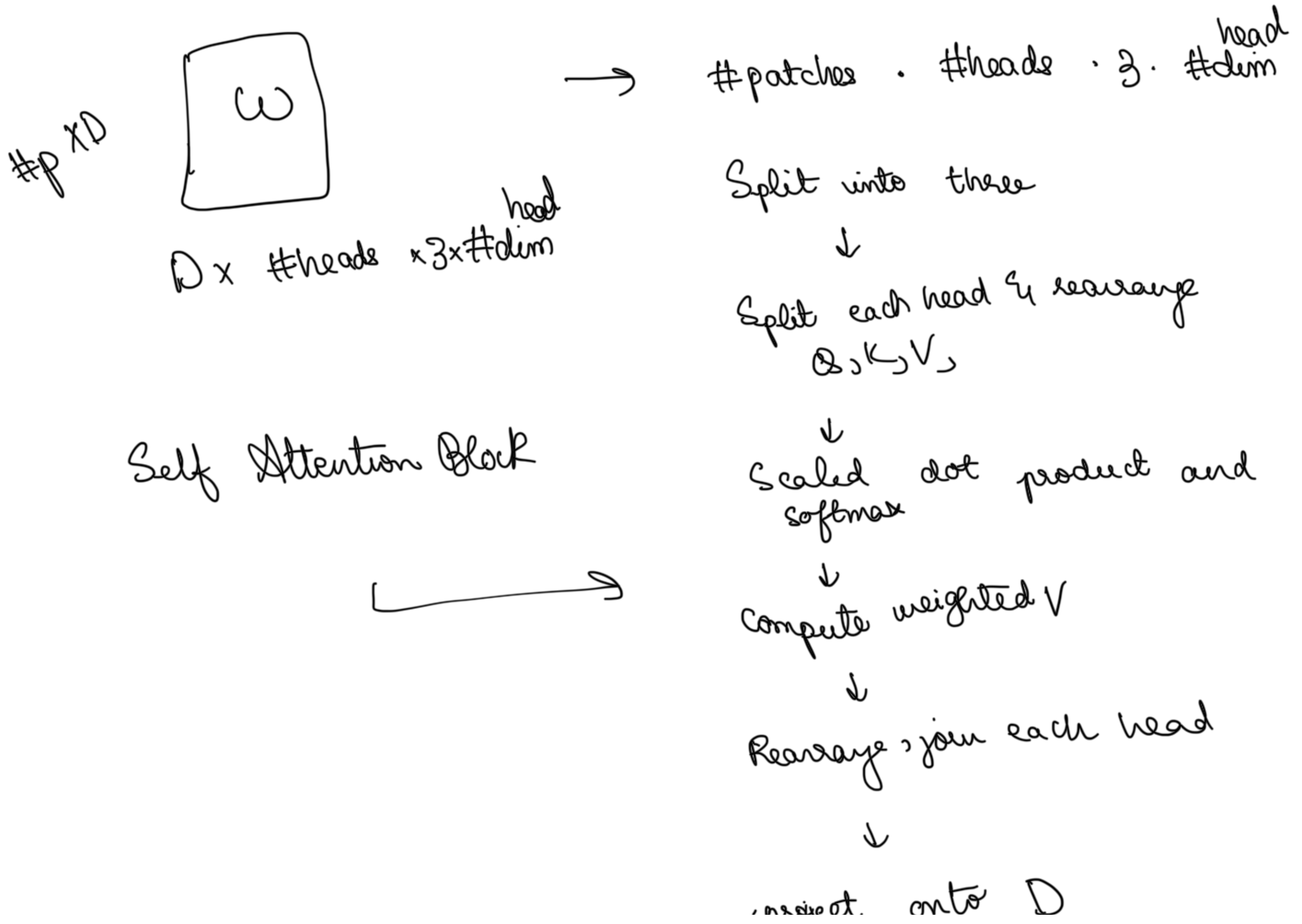
$4 \times \text{head dim}$

context  
representation  
of patch  $i$

→ Relevance for multiple factors  
→ multiple heads

$4 \times \# \text{heads} \times \text{head dim}$  →  $f_c$  layer →  $4 \times D$   
concatenate  
 $W$

So one  $W$  is enough, split to unback into  $Q, V, K$



guyana