

Self attention

$$X \in \mathbb{R}^{N \times S \times w}$$

where N = Batch size

S = seq length

w = width of embedding dimension

$$W_{QKV} \in \mathbb{R}^{3w \times w}$$

where $3w$ is QKV rows

w is width.

$$X @ W_{QKV}^T \in \mathbb{R}^{N \times S \times 3w}$$

↳ split into

$$Q \in \mathbb{R}^{N \times S \times w}$$

$$K \in \mathbb{R}^{N \times S \times w}$$

$$V \in \mathbb{R}^{N \times S \times w}$$

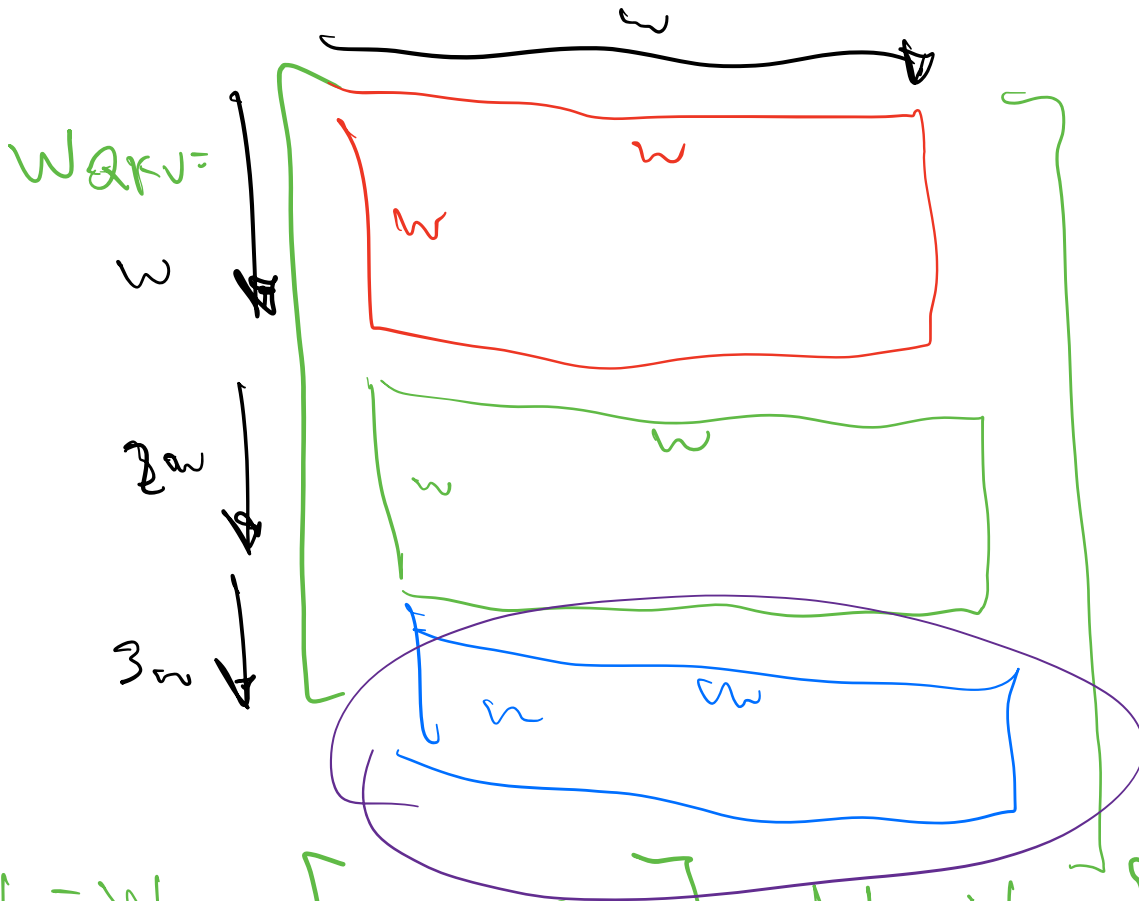
Indexing Logic :-

Just focus on (S, w) from X and $(3w, w)$

from W_{QKV}

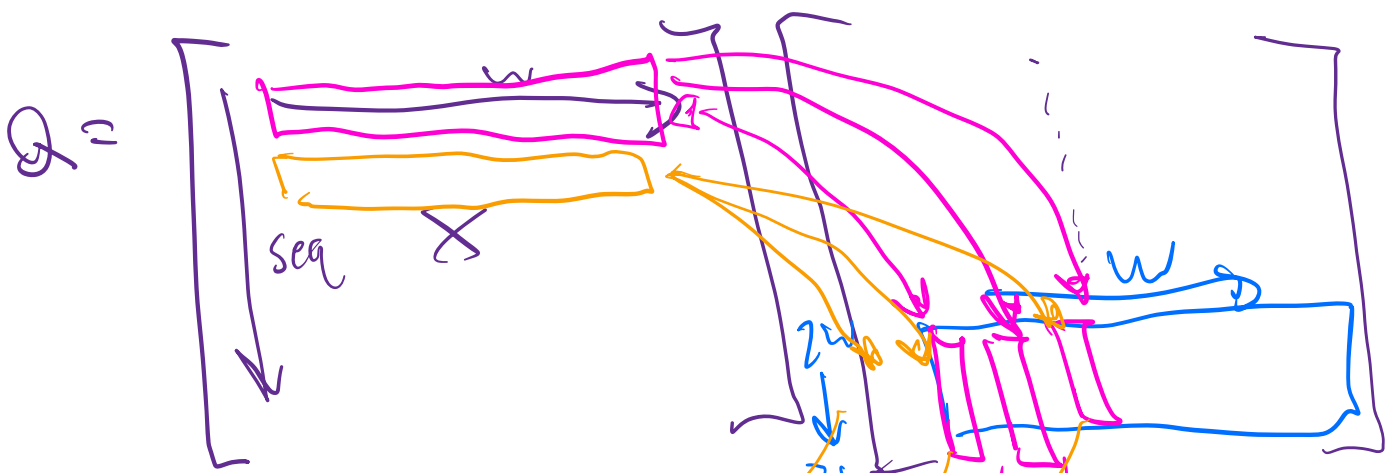
w

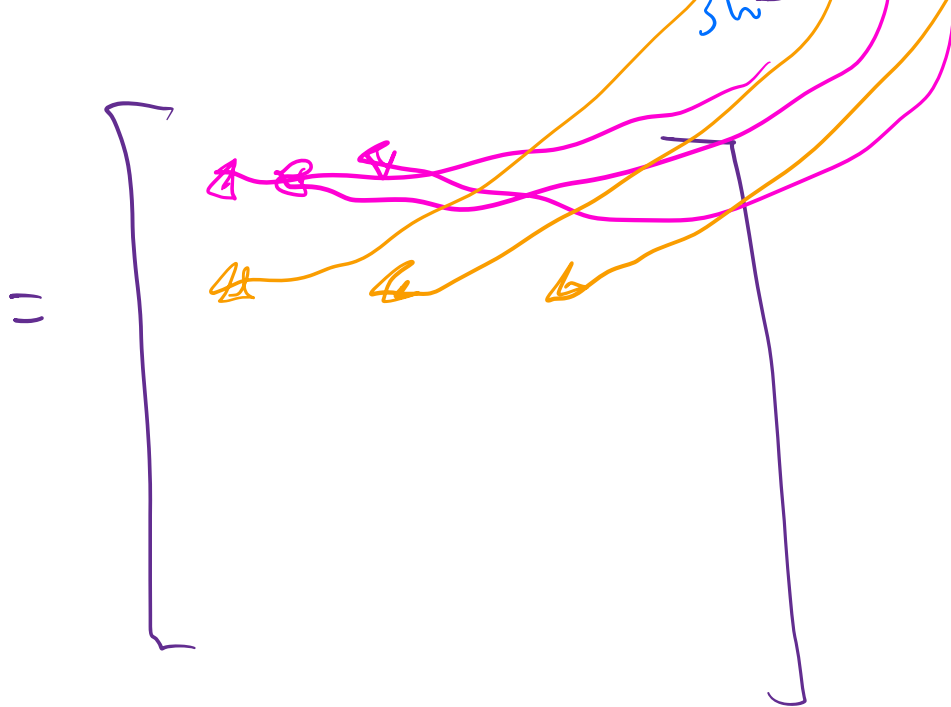
$$X = \begin{bmatrix} x_{11} & \dots & x_{1w} \\ \vdots & & \vdots \\ x_{seq1} & \dots & x_{seqw} \end{bmatrix}$$



$$W_Q = W_{QKV} [0:w, 0:w] \quad W_V = W_{QKV} [2w:, 0:w]$$

$$W_K = W_{QKV} [w:2w, 0:w]$$





$$Q_{ij} = \sum_k X_{ik} * X_{kj}$$

$$k \in [0, w-1] \text{ or } [w, 2w-1] \text{ or } [2w, 3w-1]$$

$$i \in [0, seq]$$

$$j \in [0, w-1]$$

for k in range(w):

$$Q_{ij} += X[i, k] * w[k, j]$$

$$K_{ij} += X[i, k] * w[w+k, j]$$

$$V_{ij} += X[i, k] * w[2w+k, j]$$

for $i \rightarrow \text{seq}:$

for $j \rightarrow w:$

for $k \rightarrow m:$

$$Q_{ij} += x[i,k] w[k,j]$$

$$K_{ij} += x[i,k] w[w+k,j]$$

$$V_{ij} += x[i,k] w[2w+k,j]$$

TVM Notation

for i, j, k in T.grid(seq, w, w):

$v_i, v_j, v_k = \text{T.axis.remap("sss", [seq, w, w])}$

$$Q[v_i, v_j] += x[v_i, v_k] * w[v_k, v_j]$$

$$K[v_i, v_j] += x[v_i, v_k] * w[w+v_k, v_j]$$

$$V[v_i, v_j] += x[v_i, v_k] * w[2w+v_k, v_j]$$