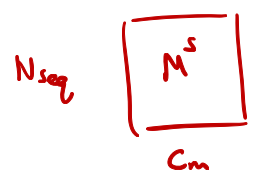
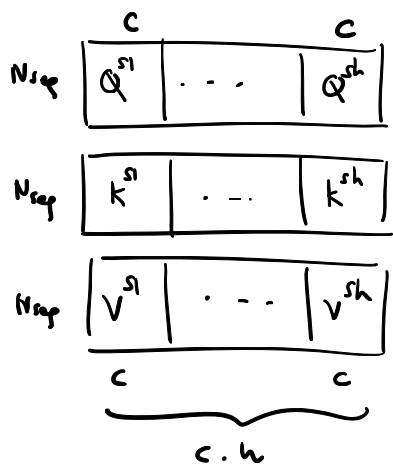


row-wise attn

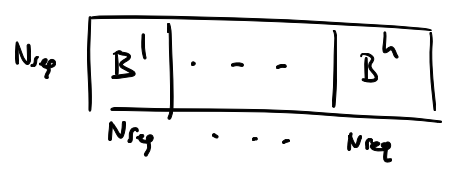
$$\forall s = 1 \dots N_{\text{clust}}$$



joint-multihead attn



$C = 32$



Shared among all s .

$$A^s = \text{softmax}_{\max} \left(\begin{array}{c|c|c} \frac{Q^{s1} K^{s1T}}{\epsilon} & \dots & \frac{Q^{sh} K^{shT}}{\epsilon} \\ \hline + B^1 & \dots & + B^h \end{array} \right)_{N_{\text{seq}}}$$

- $Q^{s1} : N_{\text{seq}} \times C$
- $K^{s1T} : C \times N_{\text{seq}}$
- $QK^T : N_{\text{seq}} \times N_{\text{seq}}$

$$G^s = \begin{array}{c} N_{\text{seq}} \\ \hline G^{s1} \dots G^{sh} \\ \hline c \end{array}$$

$$O^s = \begin{array}{c} N_{\text{seq}} \\ \hline G^{s1} \odot (A^{s1} \cdot V^{s1}) \dots G^{sh} \odot (A^{sh} \cdot V^{sh}) \\ \hline O^{s1} \dots O^{sh} \end{array}$$

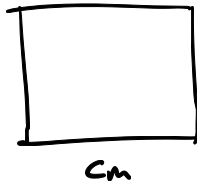
$$G^{s1} : N_{\text{seq}} \times c$$

$$A^{s1} \quad V^{s1} : N_{\text{seq}} \times N_{\text{seq}} \quad N_{\text{seq}} \times c$$

$$\underbrace{\quad \quad \quad}_{N_{\text{seq}} \times c} \quad \text{c.h}$$

$$M^s = \text{Update}$$

N_{seq}

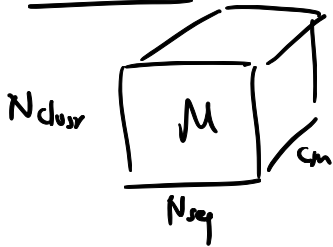


C_m

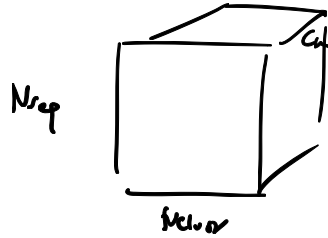
linear

$C \cdot h$
 \downarrow
 C_m

Col-wise

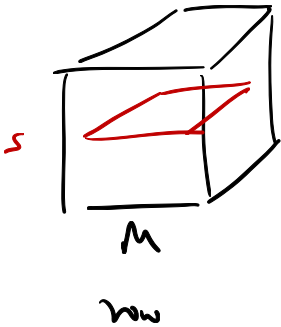


transpose

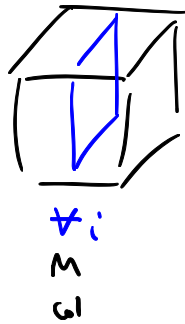


Now apply row-wise attn, without the bias matrix B^i

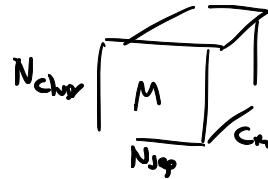
$\forall s$



→

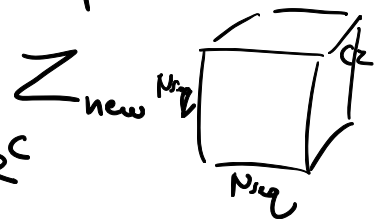


→

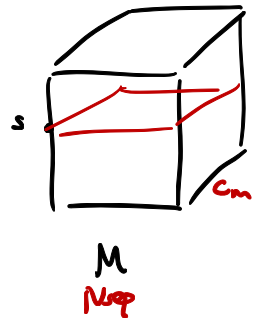


↓

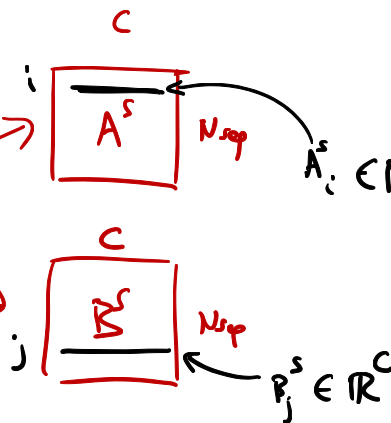
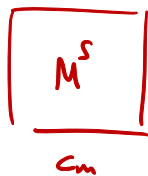
Outer product mean



N_{cls}



N_{seq}

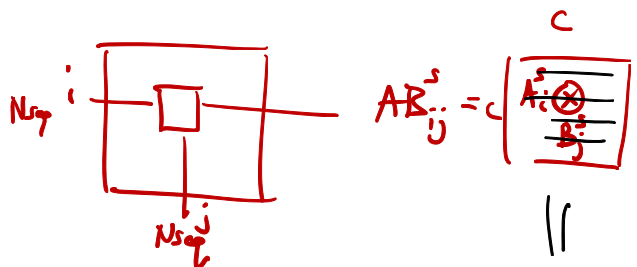


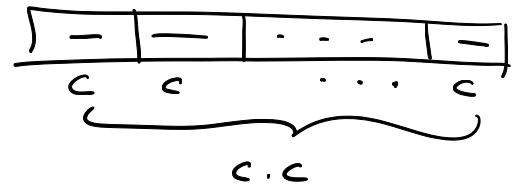
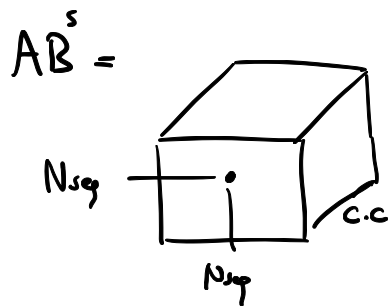
$$AB^s = A^s \otimes B^s$$

4D tensor

$$N_{seq} \times N_{seq} \times C \times C$$

↓





$$0 = \frac{1}{N_{chur}} \sum_{s=1}^{N_{chur}} AB^s$$

$$N_{seq} \times N_{seq} \times (c.c)$$

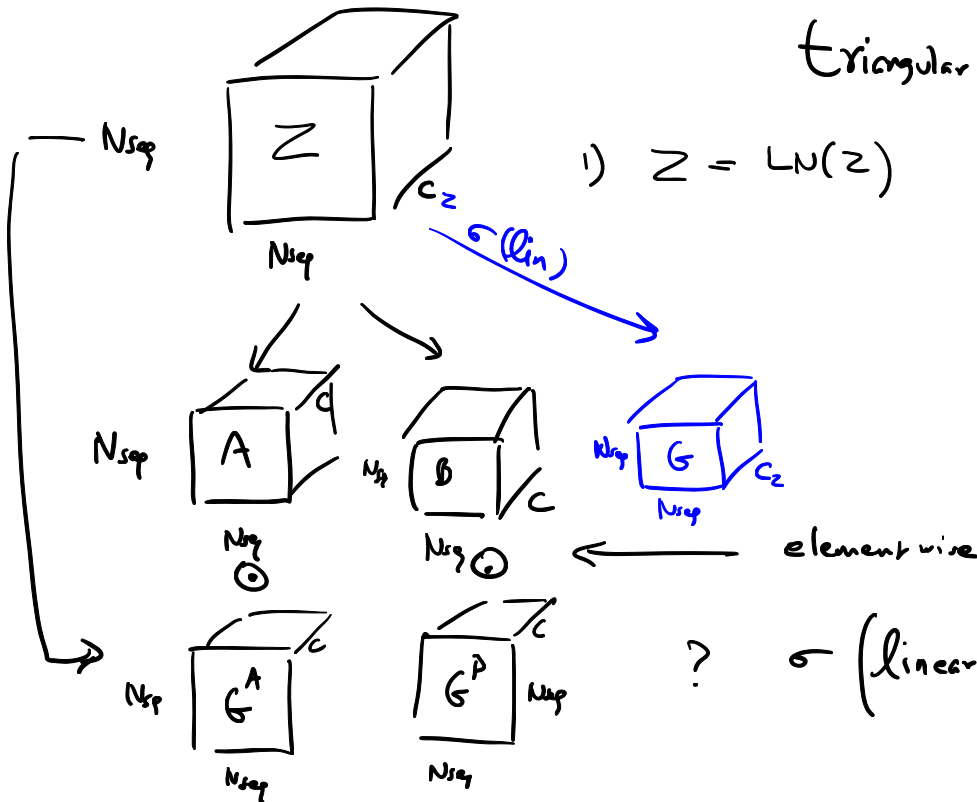


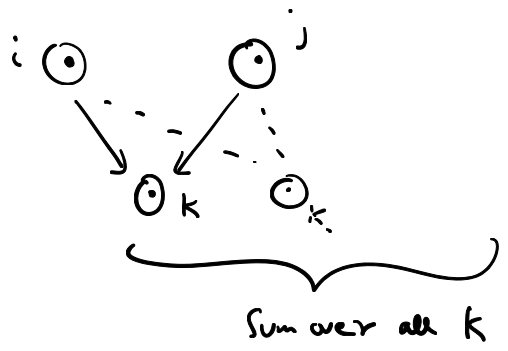
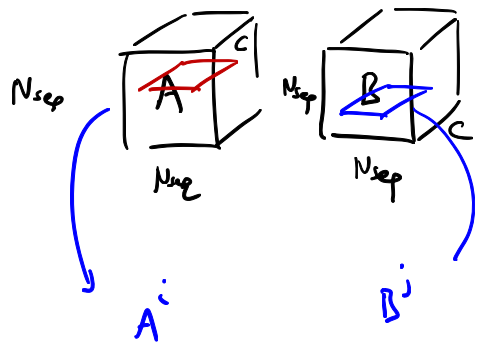
$$Z_{new} = (N_{seq} \times N_{seq} \times c_s)$$

$$Z = Z + Z_{new}$$

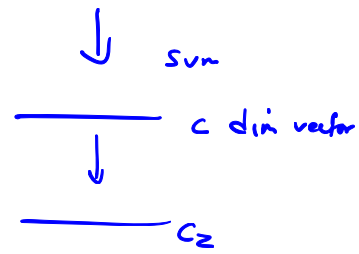
Operations on pair-sep ($c=128$) ($c_s=128$)

Triangular multiplicative update (out)

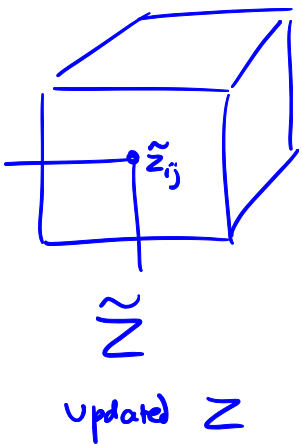




$$\underbrace{N_{sep} \times c \text{ matrix}}_{A^i} \odot \underbrace{N_{sep} \times c \text{ matrix}}_{B^j} = AB^{ij} \quad N_{sep} \times c \text{ matrix}$$

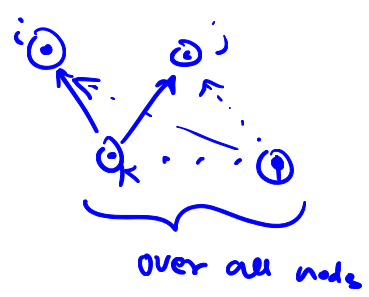
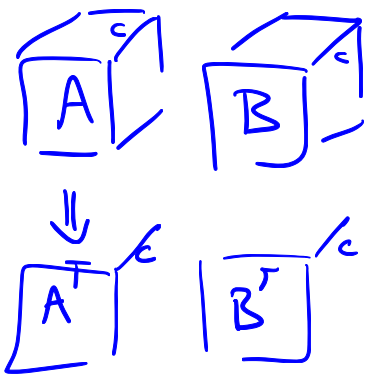


$$\tilde{z}_{ij} = \text{---} c_z$$

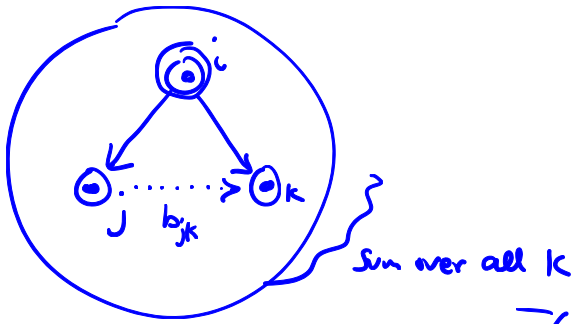


$$\tilde{Z} = \tilde{Z} \odot G \quad \text{gate} \quad N_{sep} \times c_z \text{ tensor } G$$

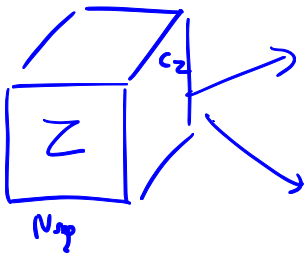
Multiplicative triangular update (incoming)



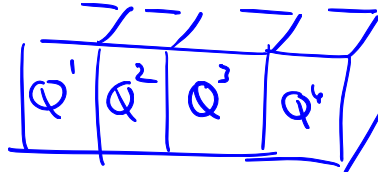
Triangular Attn (starting)



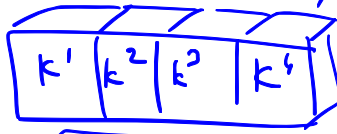
N_{seq}



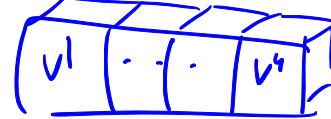
Q



K



V

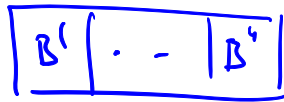


$h = 4$

$c = 32$

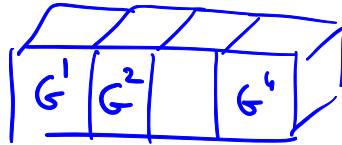
$(N_{seq} \times N_{seq} \times c)^{x h}$ ← # of heads

B



bias $(N_{seq} \times N_{seq})$

G



$(N_{seq} \times N_{seq} \times c)^{x h}$

$\forall h = 1 \dots 4$

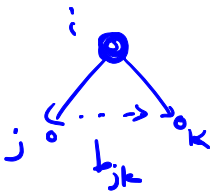
$$a_{ijk} = \text{softmax} \left(\frac{\vec{q}_{ij}^T \vec{k}_{ik}}{\sqrt{c}} + b_{jk} \right)$$

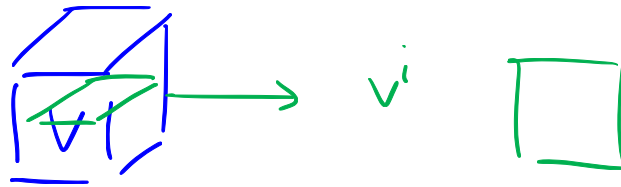
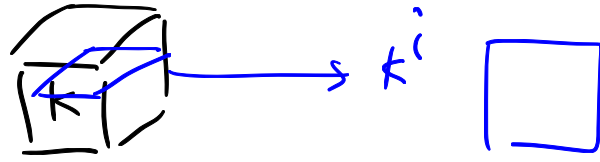
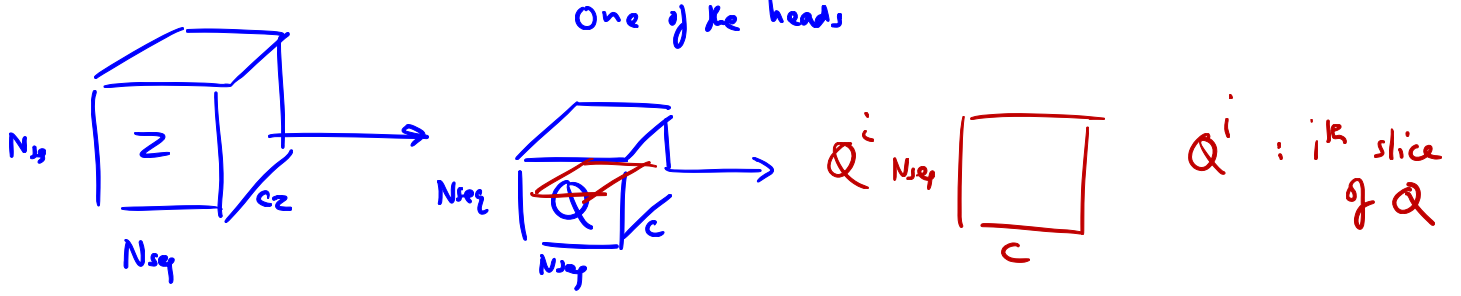
Attn between i & j wr.t k

$$\underline{O}_{ij} = \frac{\vec{g}_{ij}}{c} \odot \sum_k \underset{\text{factor}}{a_{ijk}} \cdot \frac{\vec{v}_{ik}}{c}$$

$$\tilde{Z}_{ij} = \text{lin}(\text{concat} (O_{ij}^h))$$

$c \cdot h$
↓
 c_2



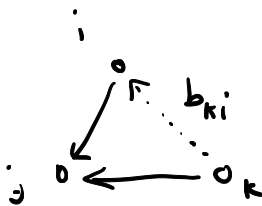


$$A^i = \left(\frac{Q^i K^i T}{\sqrt{c}} + B^T \right)$$

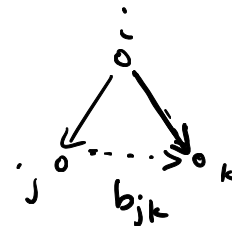
check this

$$O^i = \text{GO}(A^i \cdot V^i)$$

triangular attn (ending)



starting

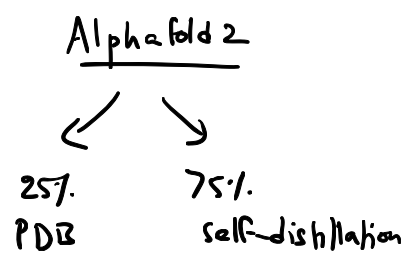
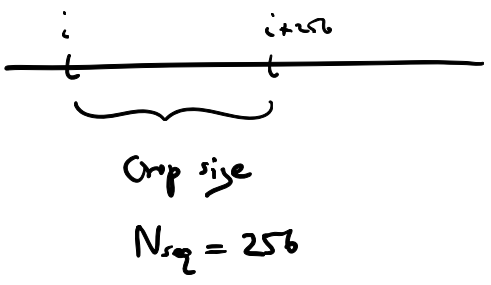
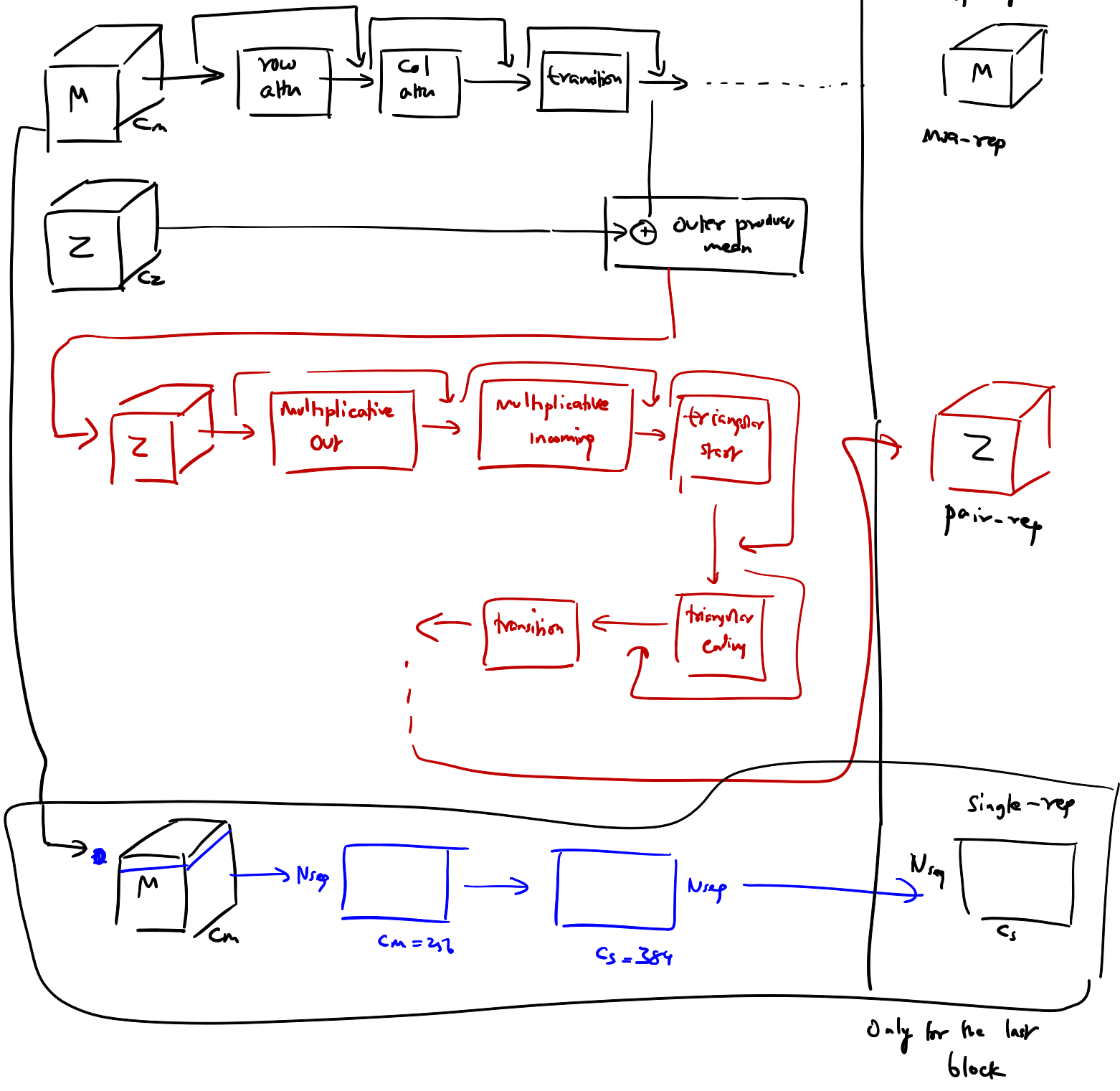


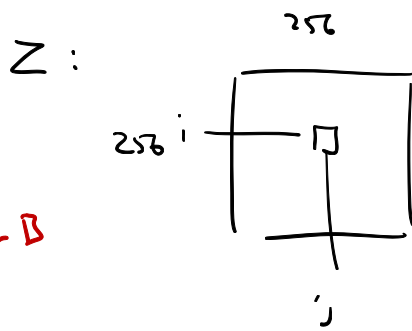
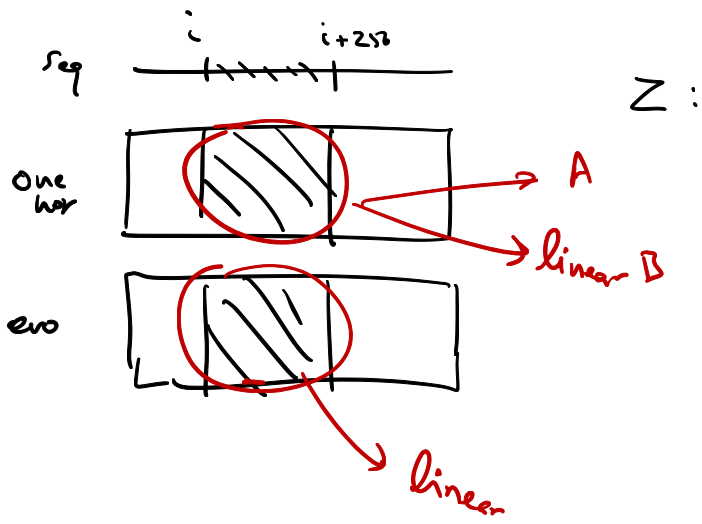
$$a_{ijk} = s \left(\frac{Q_{ij}^T K_{kj}}{\sqrt{c}} + b_{ki} \right)$$

only change

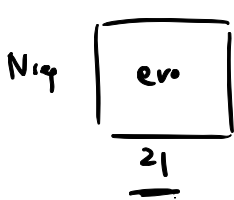
$$\text{softmax} \left(\frac{Q^i (K^i)}{\sqrt{c}} + B^T \right)$$

check this!



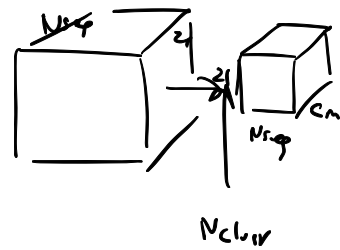


$$z_{ij} = a_i + b_j$$



Outer product

$$N_{seq} \times 21 \times 21 \rightarrow 21$$



16 {linear p_{ij}
 $N_{clust} = 16$

