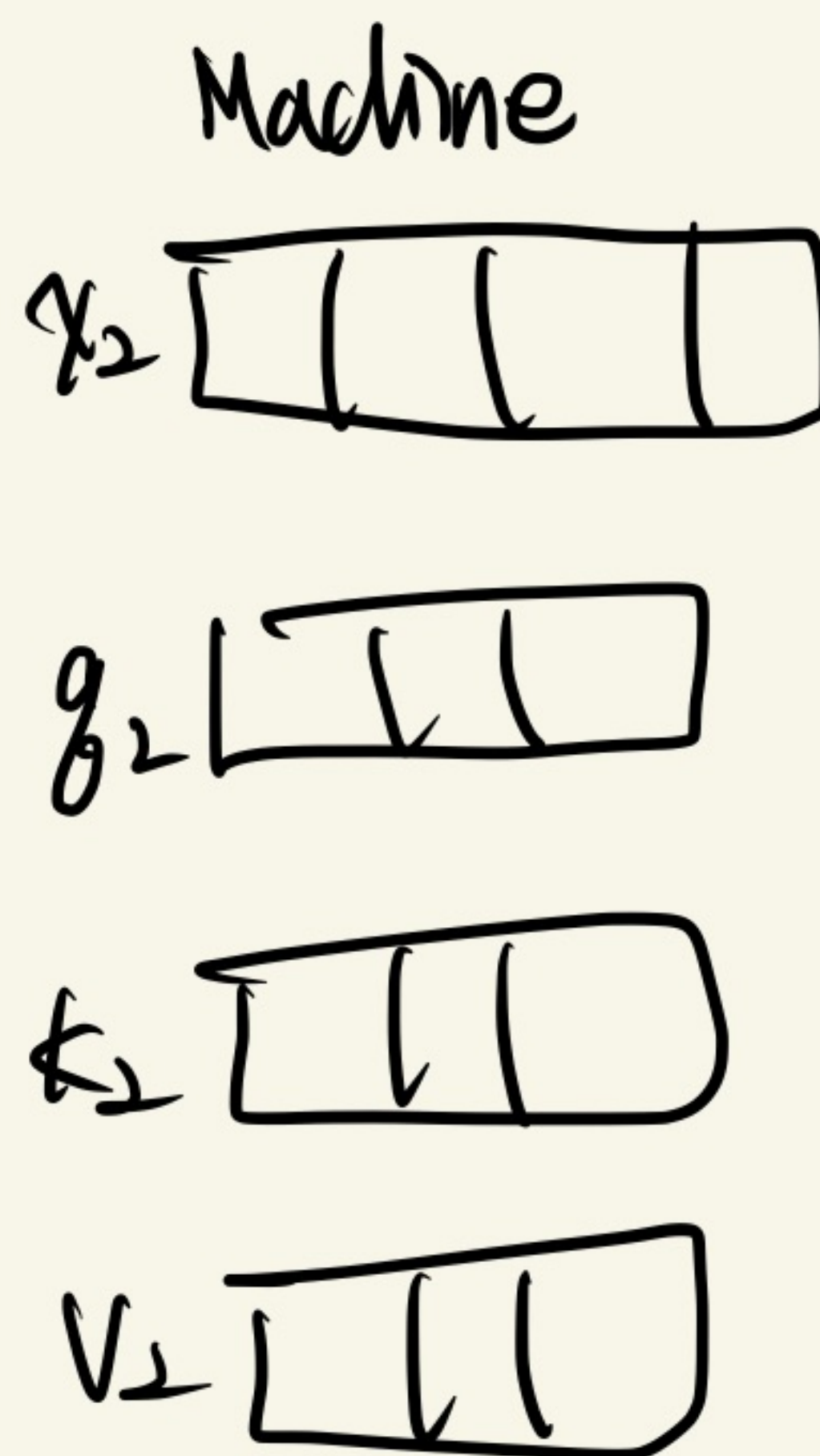
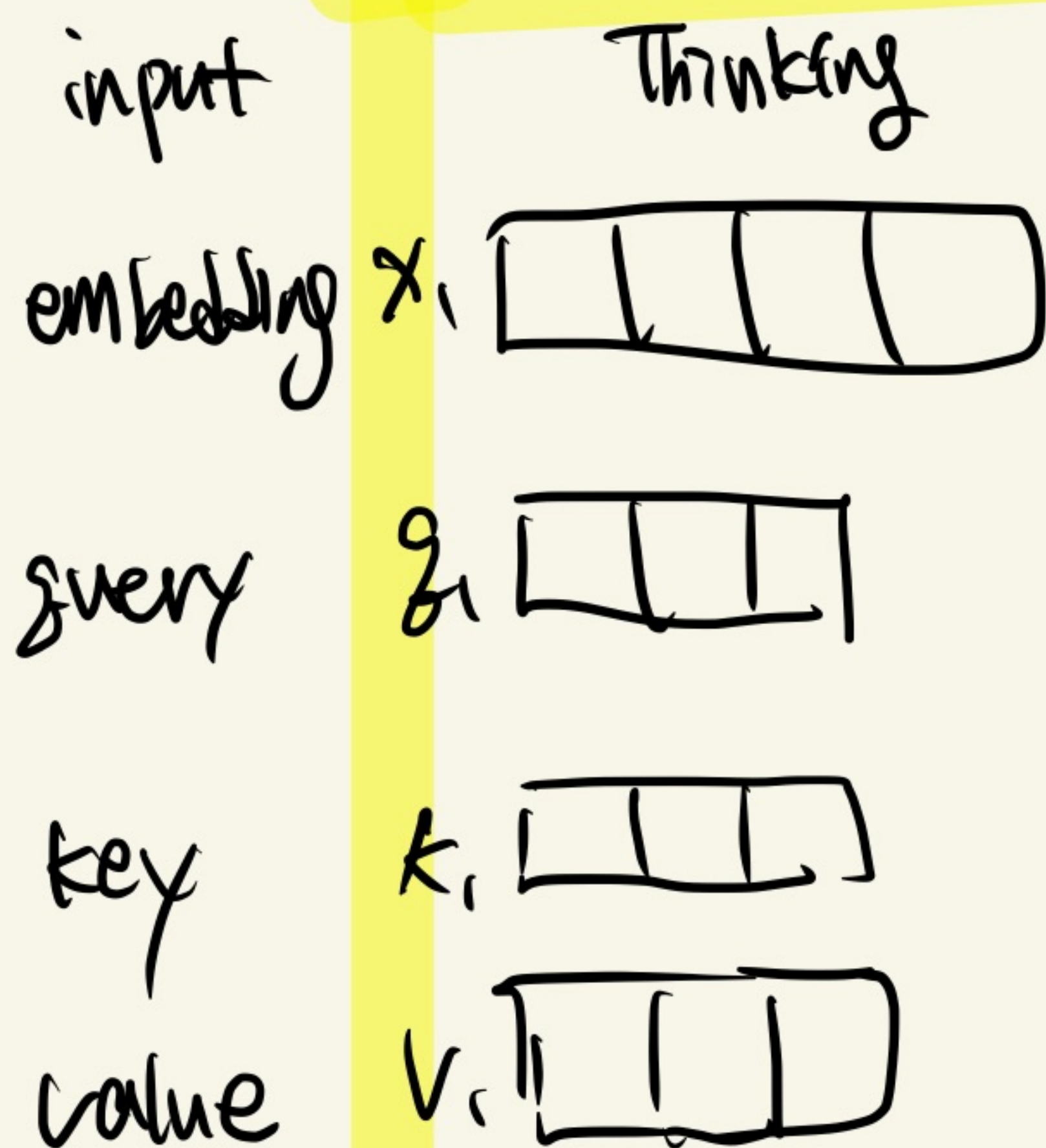


#1 $x_1 \dots x_n \rightarrow z_1 \dots z_n$ process



score $q_1 \cdot k_1 = 112$

divide by d
($\sqrt{d_k}$)

14

Softmax

0.99

$q_1 \cdot k_2 = 96$

12

0.12

value v_1

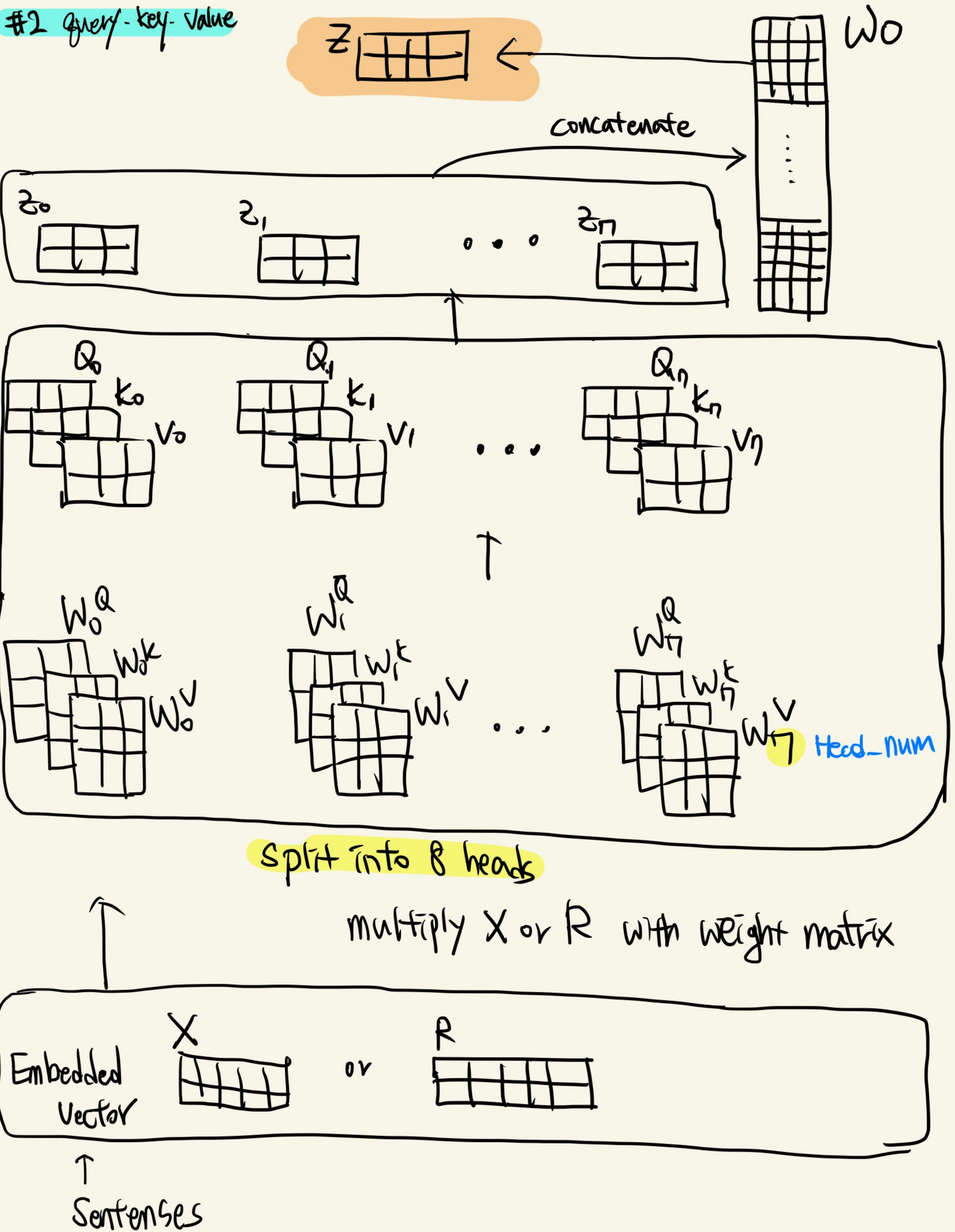
v_2

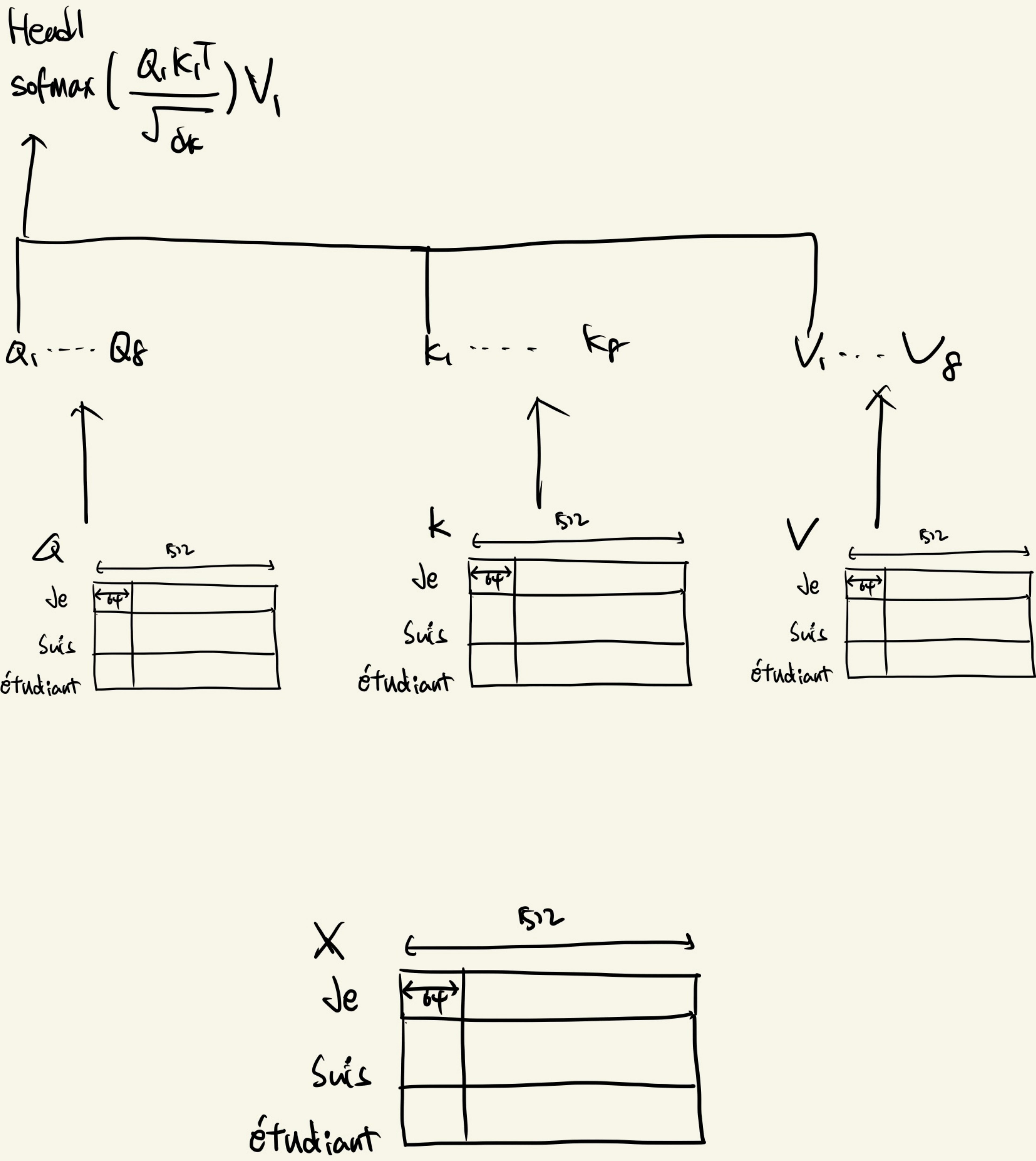
Sum z_1

z_2

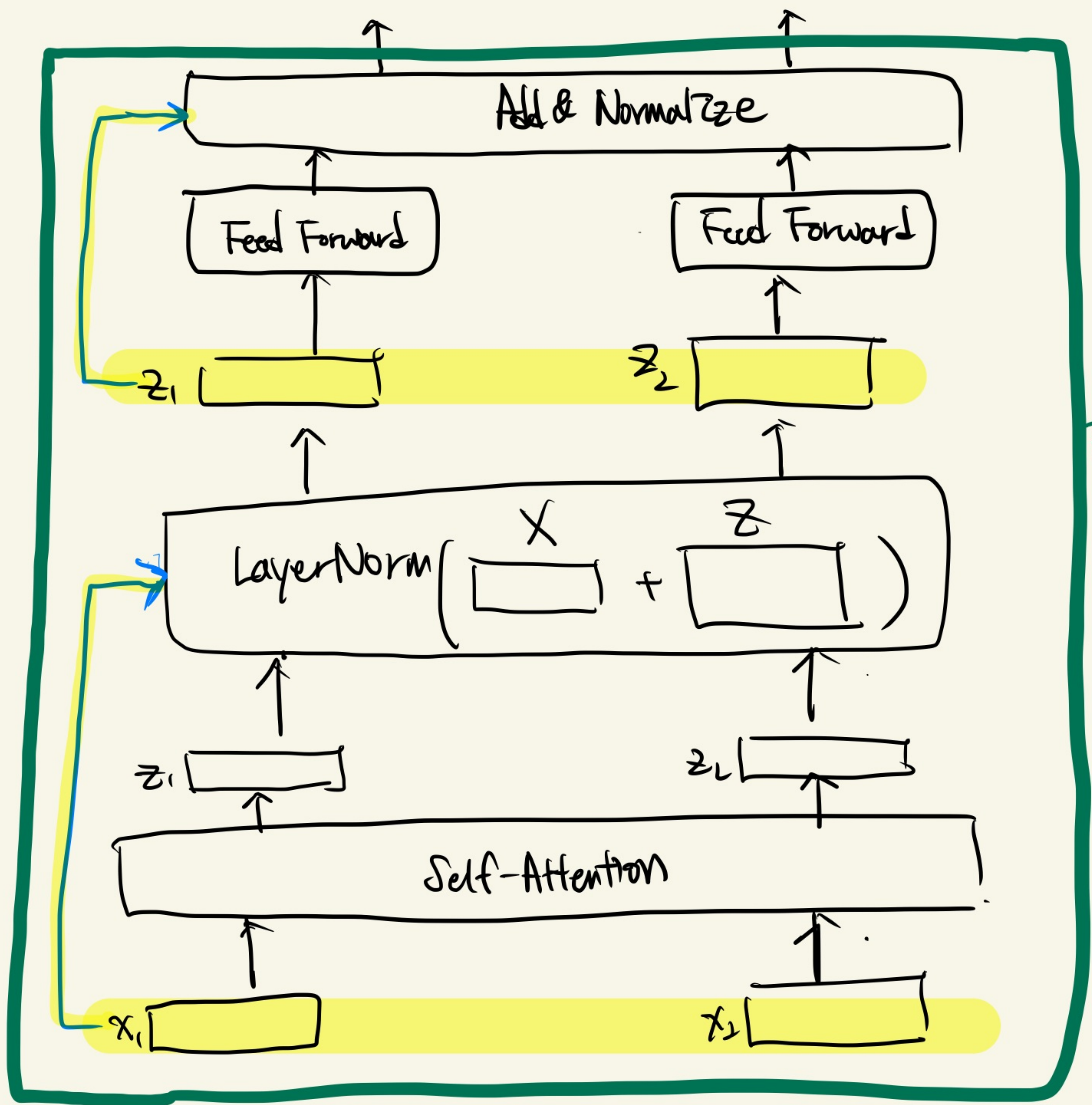
$$Z = \text{Softmax}\left(\frac{Q \cdot K^T}{\sqrt{d_k}}\right) V$$

#2 query-key-value





#3 Encoder

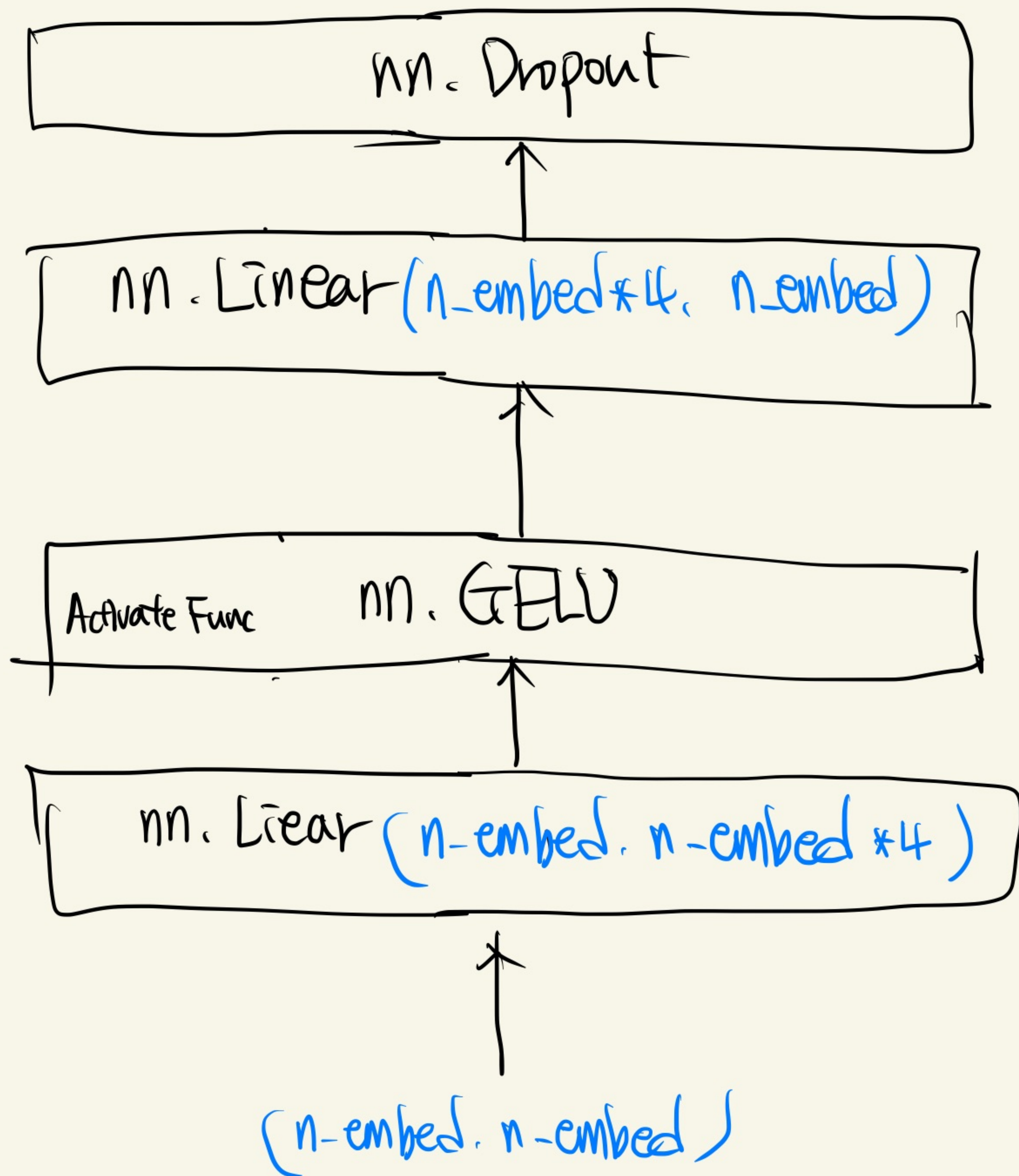


Encoder
#1

Positional Encoding \oplus
 x_1
 Thinking

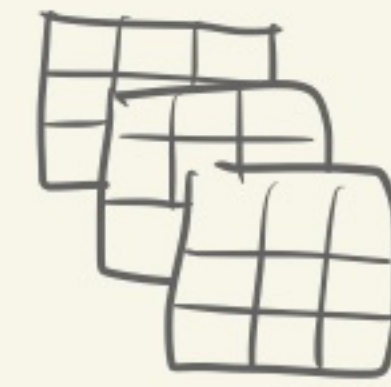
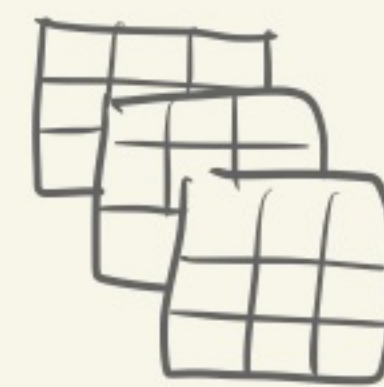
Positional Encoding \oplus
 x_2
 Machines.

#4. Feed Forward .



Entire block

$K_{encoder}$ $V_{encoder}$



Encoder #2

softmax

Linear (Logit)

Decoder #2

Add & Normalize

Feed Forward

Feed Forward

Add & Normalize

Encoder-Decoder Attention

Add & Normalize

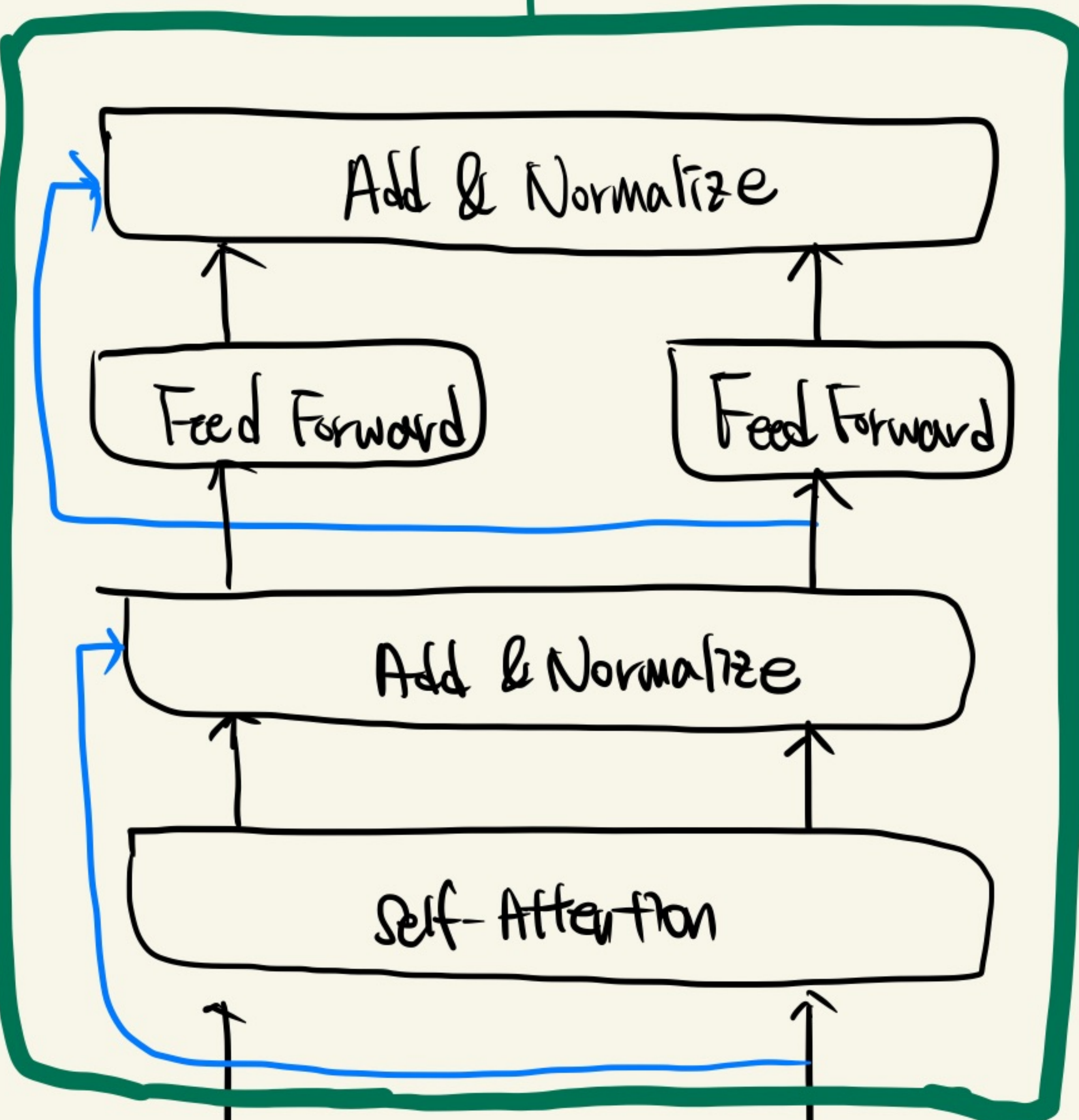
self-Attention

Decoder #1



* end of sentence η_1, η_2
* earlier position was allowed (-inf)

Encoder #1

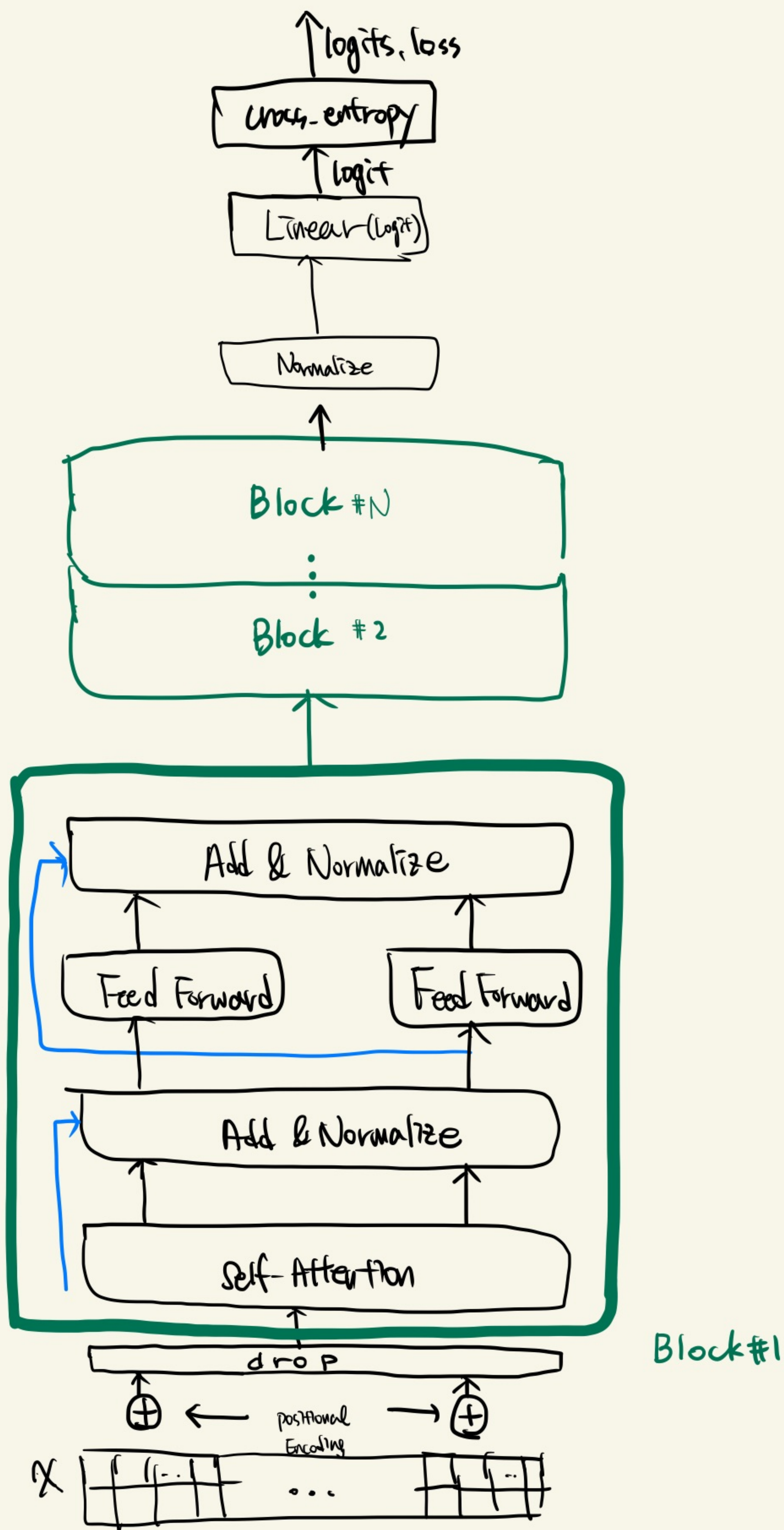


x_1 [[(-1)]]
Thinking

positional
Encoding

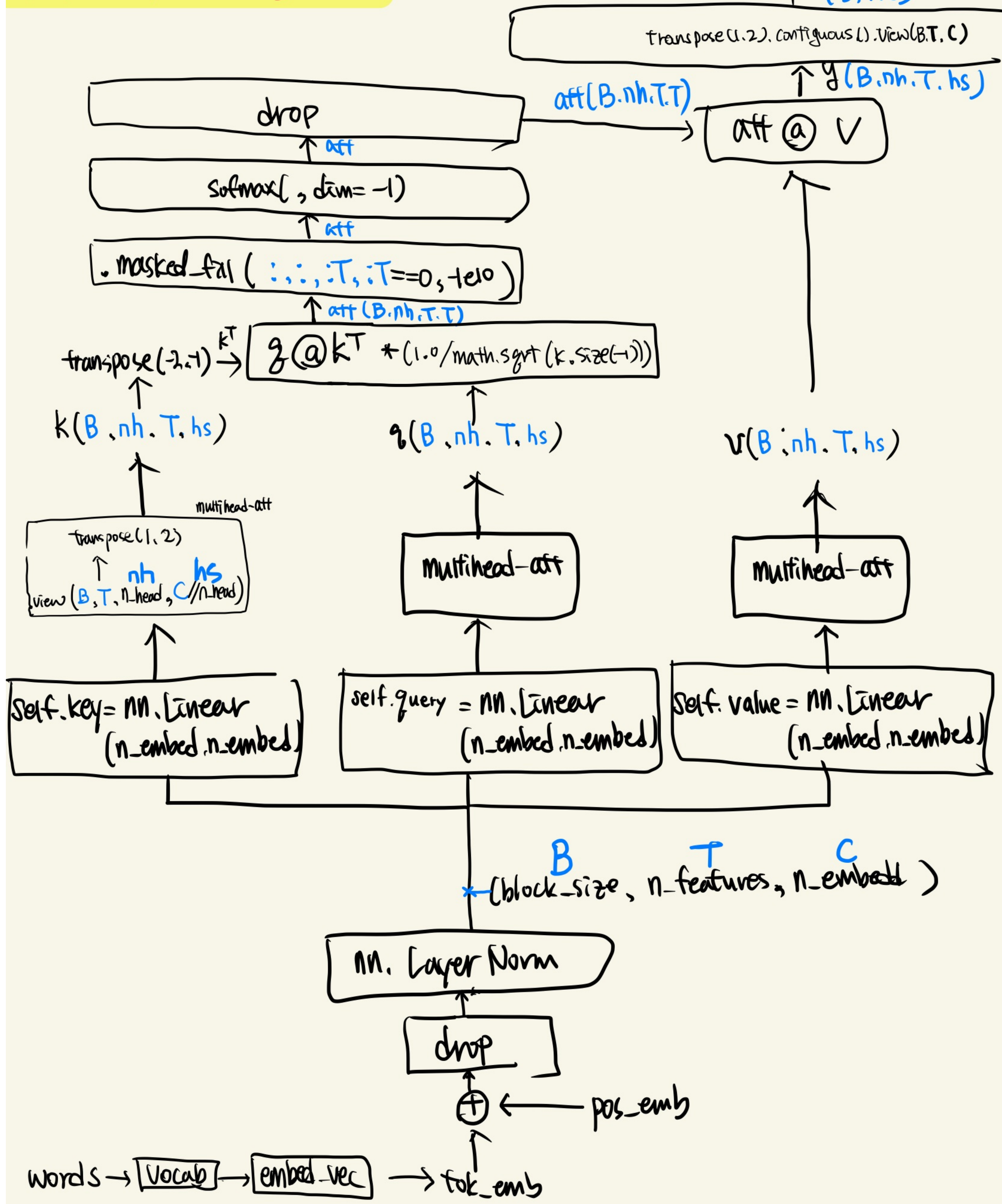
x_2 [[(-1)]]
Machine

#4 A5 GPT - decoder concept - unidirectional .

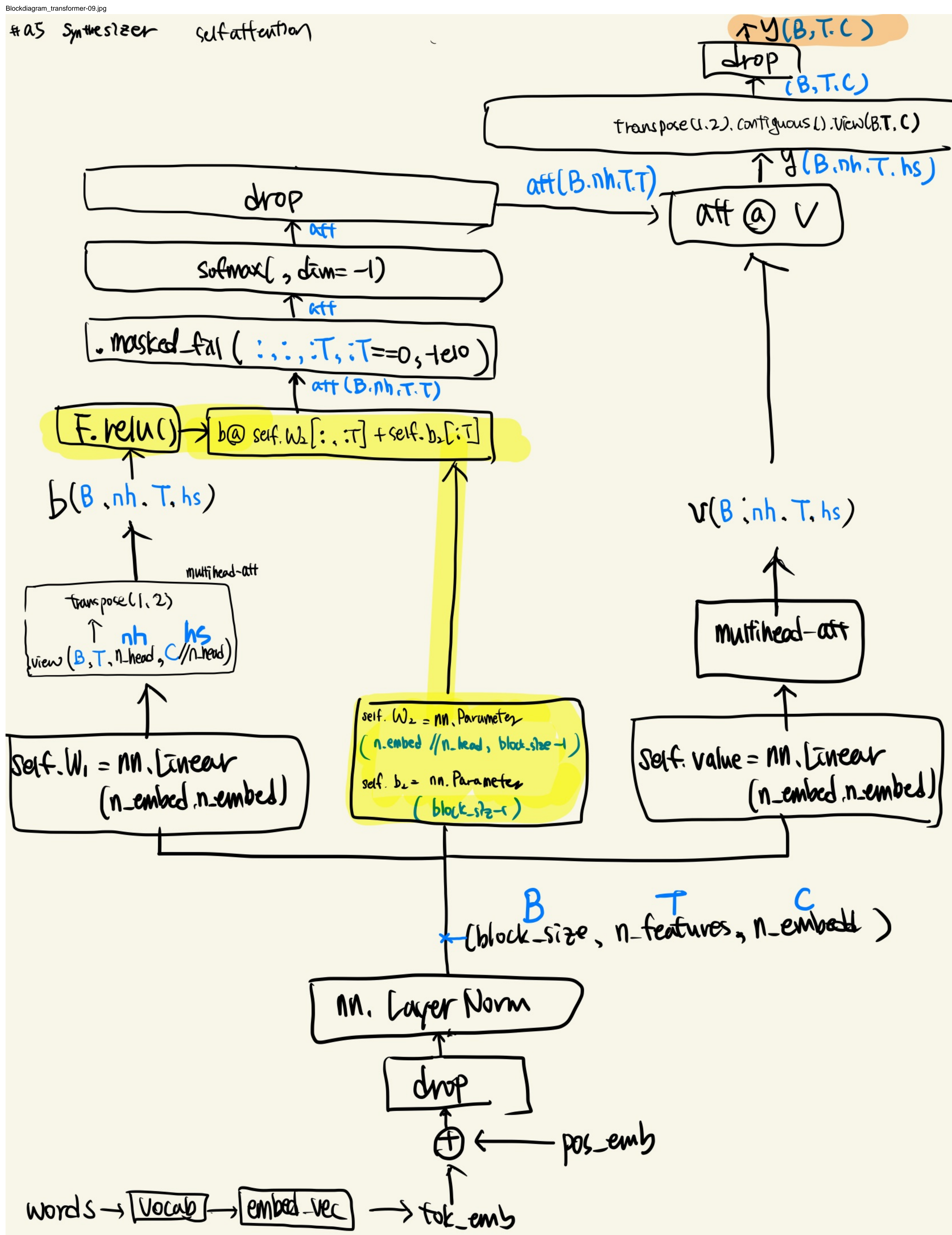


#a5 Casual selfattention

nn.Linear(bias=True) default



a5 Synthesizer self attention



a5 Synthesizer self-attention (equation version)

