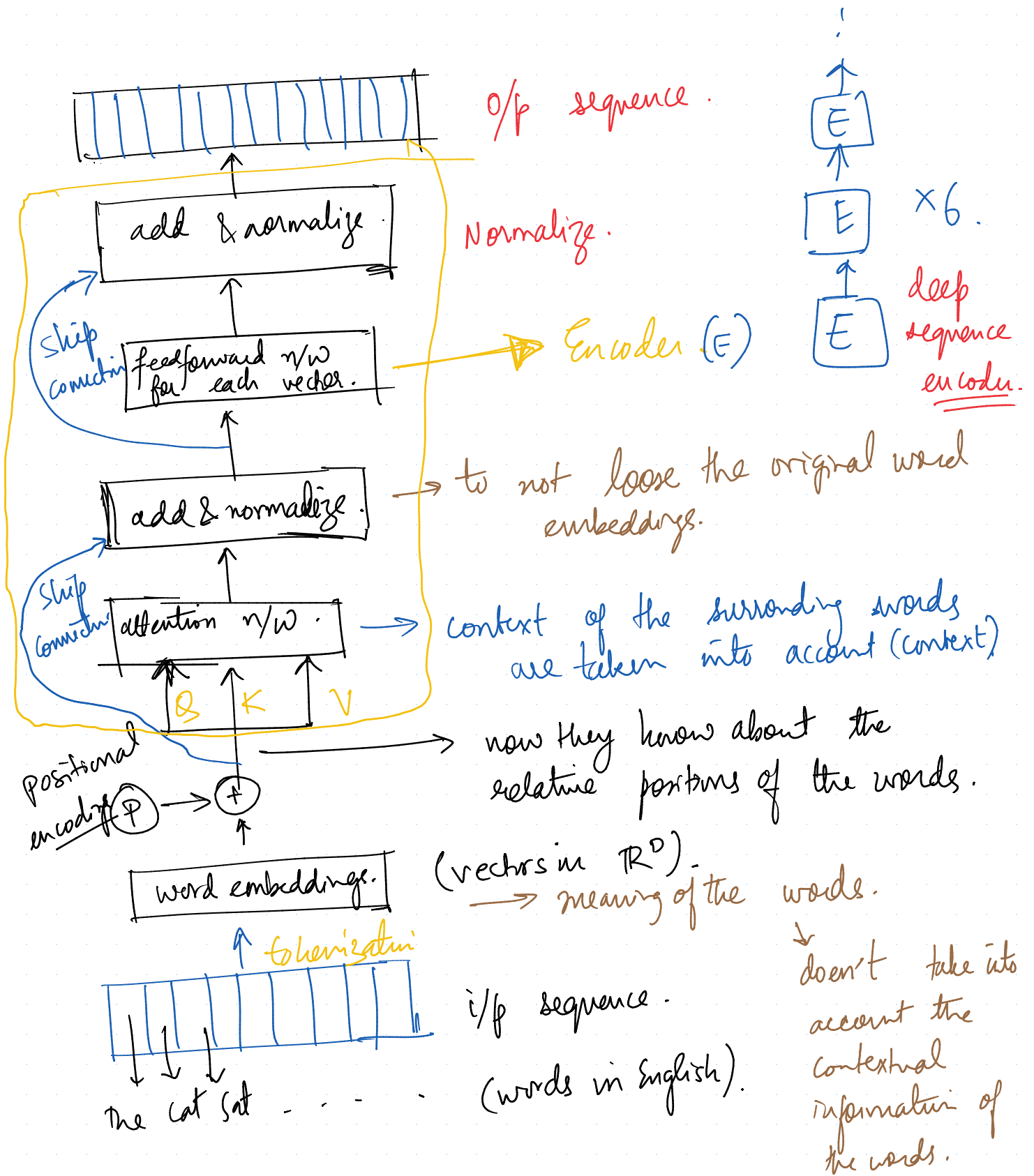


Attention - Based Sequence Encoder.



Feedforward Network.

Provide regularization or structure to the n/w.

Restrict the o/p of the n/w to be contained

to the subspace of the vectors with the n/w.
(maths concept)

tan h \rightarrow if we use \rightarrow o/p (-1 & 1)

$$v = \frac{1}{N} \sum_{n=1}^N \tilde{c}_n$$

$$p(\text{+ve sentiment}) = \frac{\exp(v \cdot e)}{1 + \exp(v \cdot e)}$$

(0-1)

Sigmoid

$$= \frac{1}{1 + \exp^{-(v \cdot e)}}$$

Add some activation like sigmoid/softmax for classification after the sequence encoder.

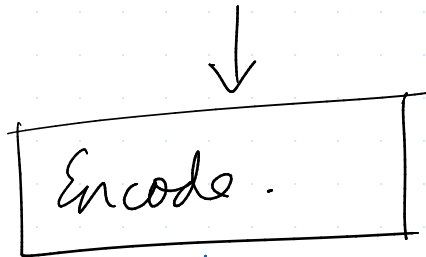
Predict the next word.

$$p(w_{n+1} = k) = \frac{\exp(v \cdot e_k)}{\sum_{j=1}^{|V|} \exp(v \cdot e_j)}$$



Predict the sequence of the next word.

English The orange cat is Nacho



Translation task.

French Le chat orange est Nacho

or UK

source text: w_1, w_2, \dots, w_n

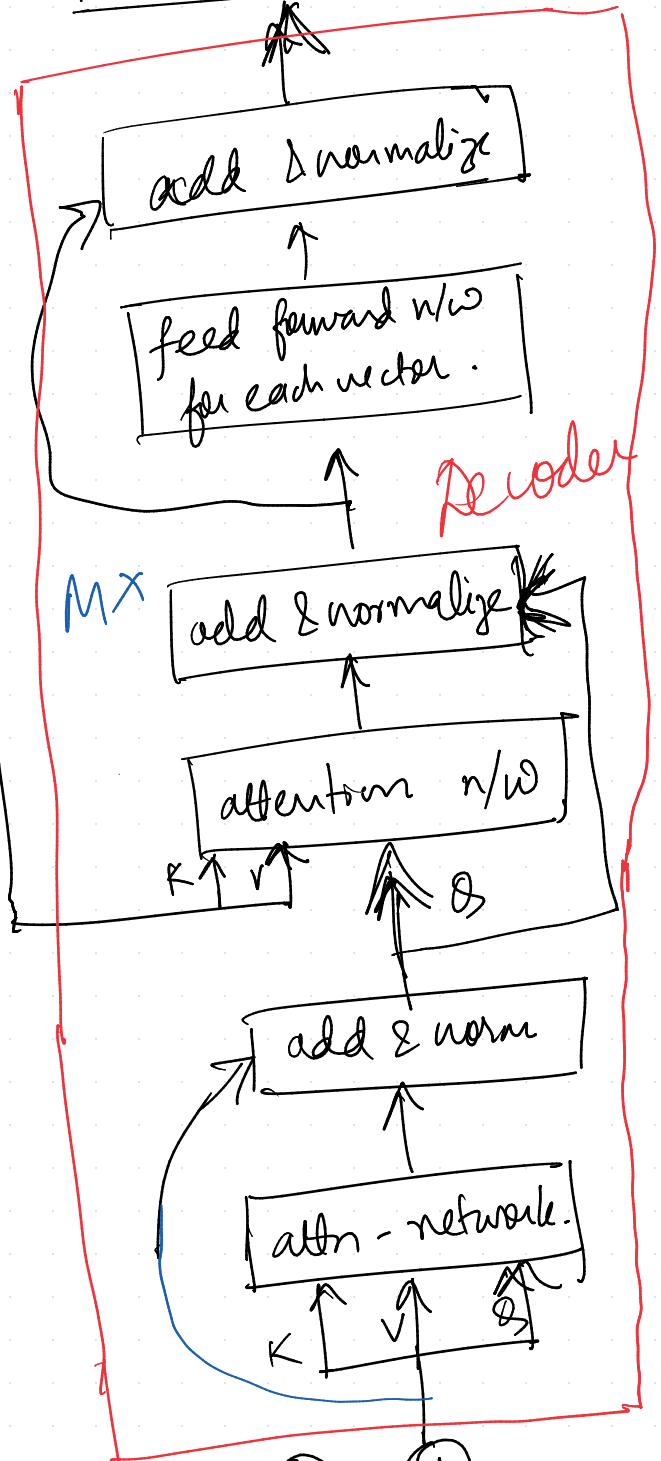
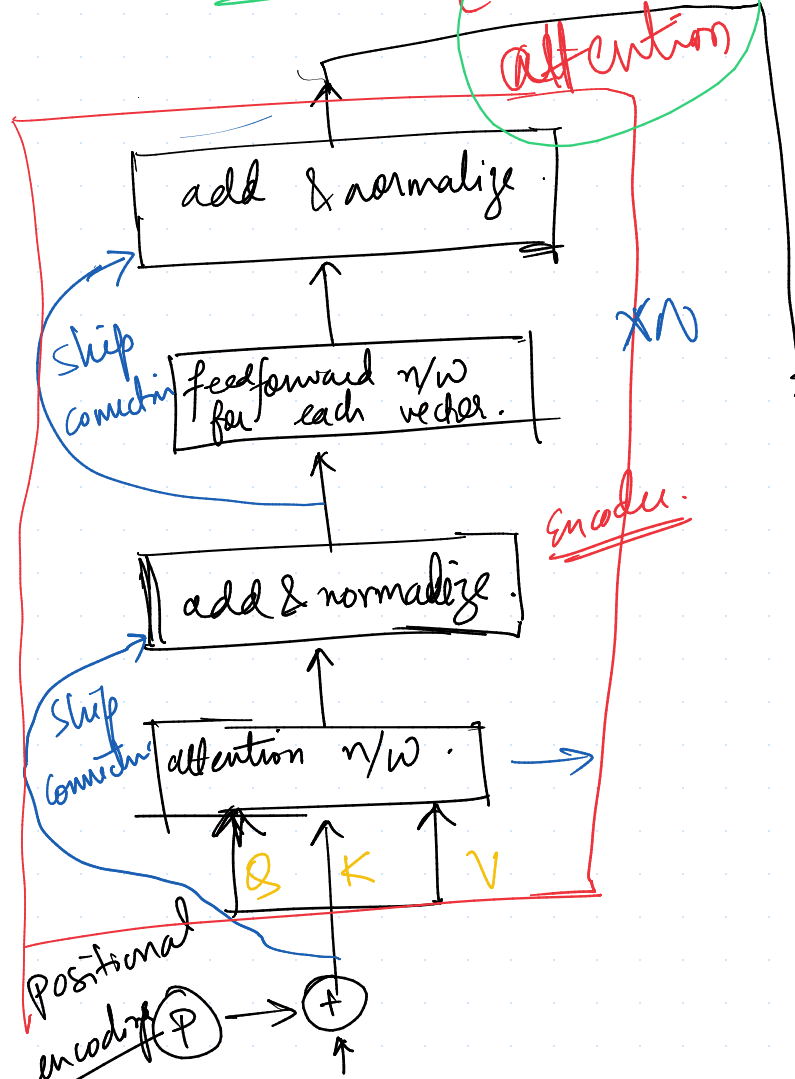
summarize. → Transformers.

Softmax predicted word.

arrange vectors.

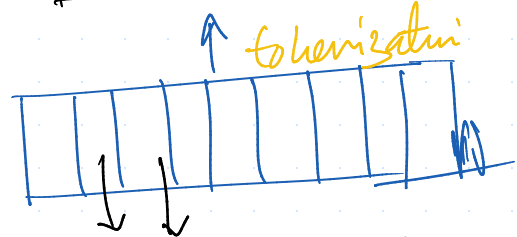
$(n \times d)$

cross attention



Positional encoding ϕ

word embeddings.

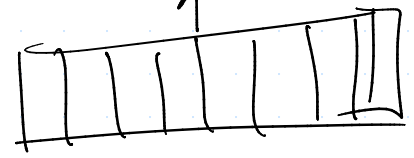


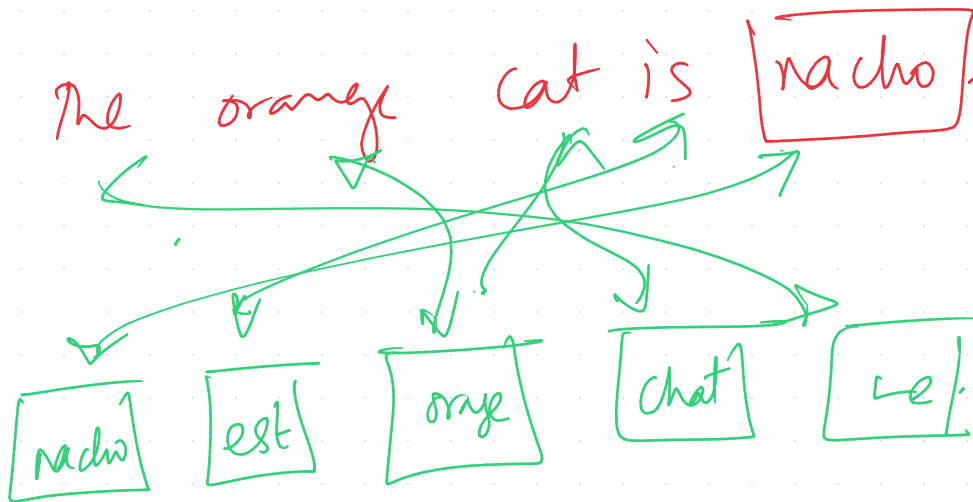
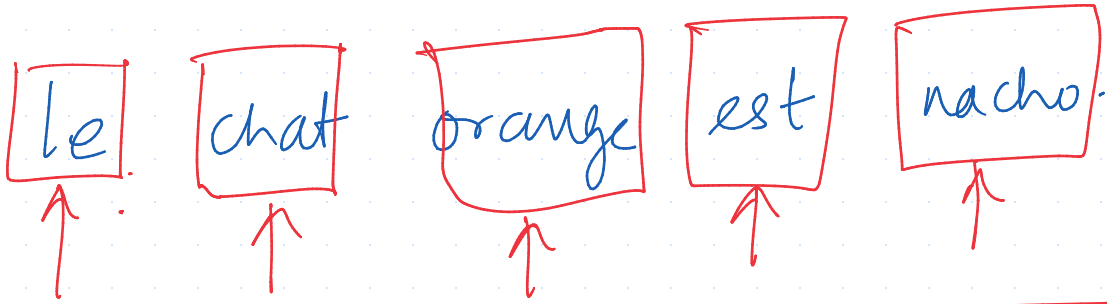
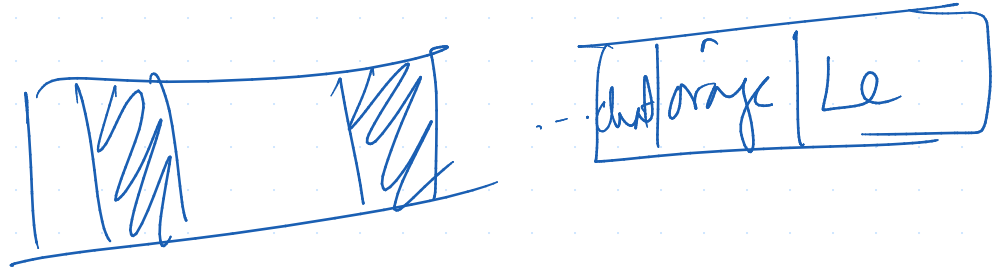
The orange cat is nacho.



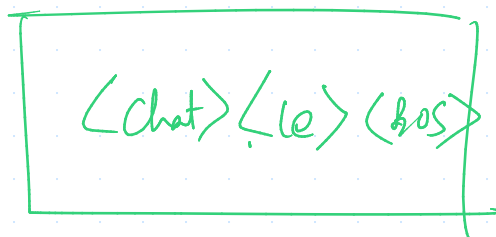
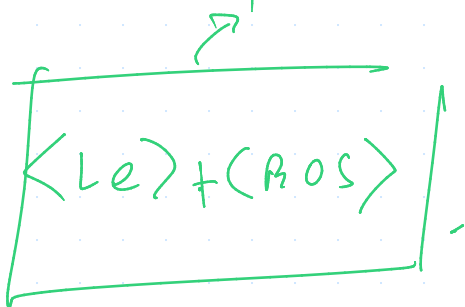
<BOS>

word embeddings.



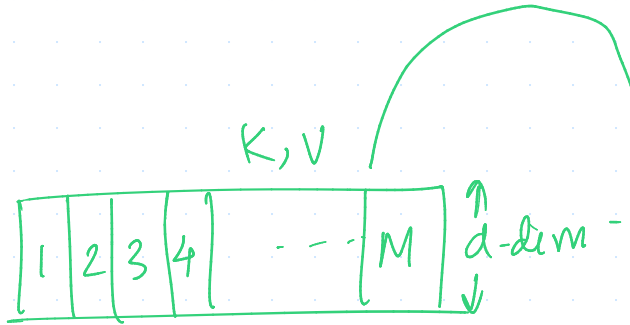


↑/p to the bottom right is shifting as we predict new words.

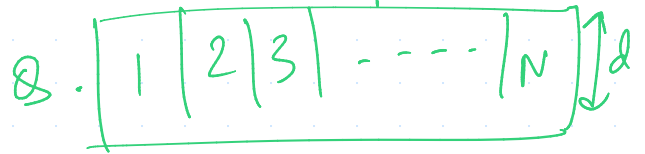
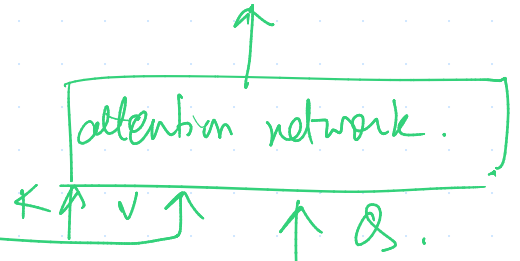
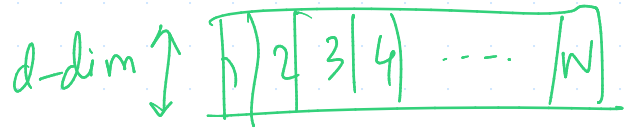


Cross attention

o/p sequence represented in terms of i/p sequence vectors



encoder top



o/p of self attention n/w applied to i/p sequence.

o/p to self-attention network applied to sequence decoder

i/p to the decoder

thus far N-words have been (8) decoded thus far → query

Cross attention :-

K, V → coming from the i/p N-words that are encoded at the top of the encoding n/w.

encoding the i/p sequence → left → Encoder.

decoding the o/p sequence → right → decoder.

• self attention \rightarrow b/w the N -words.

• Cross attention \rightarrow b/w the o/p of the encoding & the o/p of the decoding.

Multihead attention.

Output (o/p) = Attention (K, V, Q).

Q

1	2	3	4	...	N
---	---	---	---	-----	-----

 d -dim.

V

1	2	3	...	N
---	---	---	-----	-----

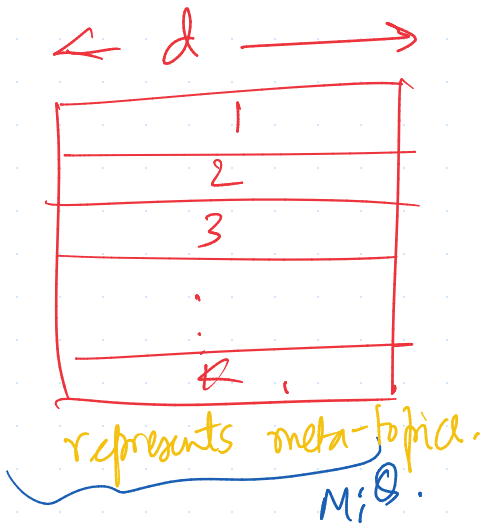
 d -dim.

K

1	2	3	...	N
---	---	---	-----	-----

 d -dim.

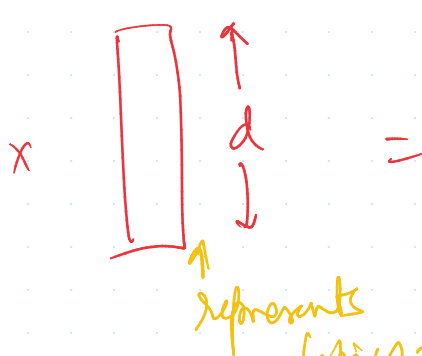
Project onto k-vectors.



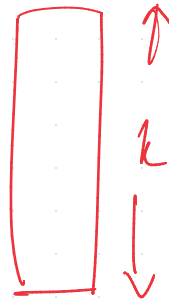
projection matrix M^Q
associated with
query vectors.

Query Matrix.

original query
 q_i



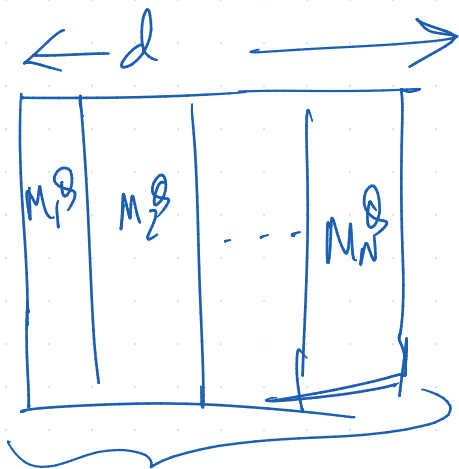
new vector



linear subspace defined by
 M^Q .

New projected vector denoted

$M^Q q_i$ applied to all
queries as $M^Q Q$.

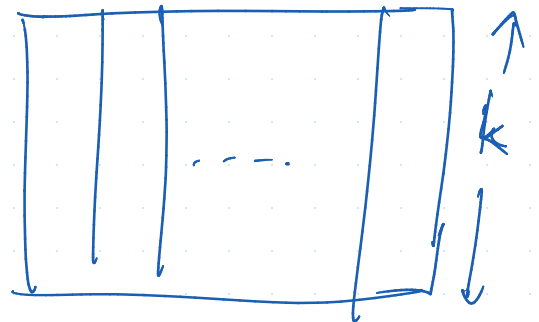


represents meta-topics
 M_i^Q .

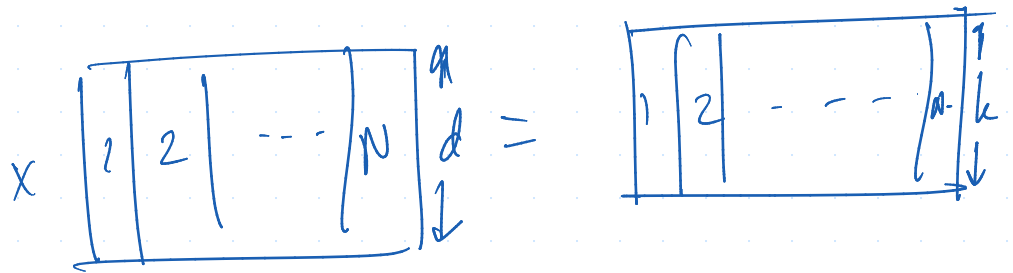
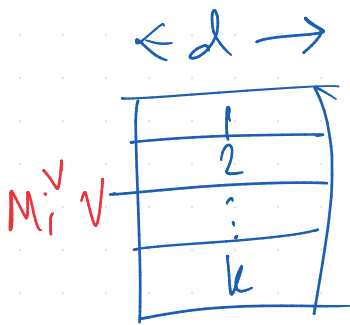
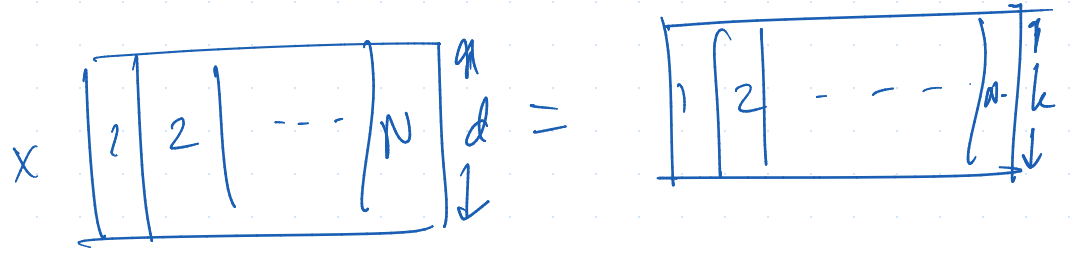
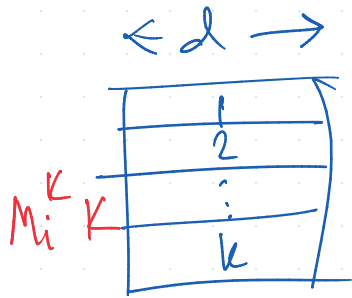
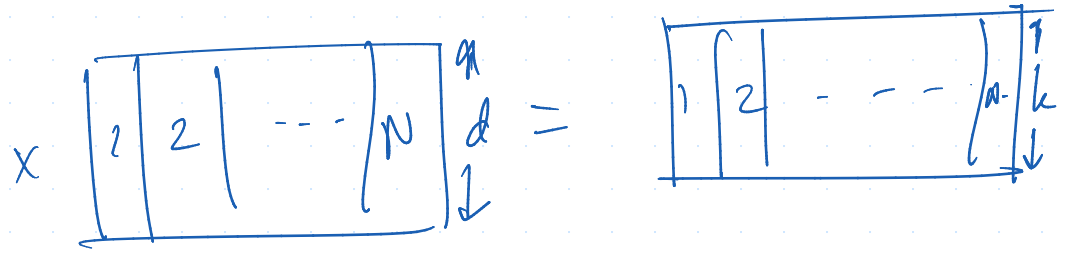
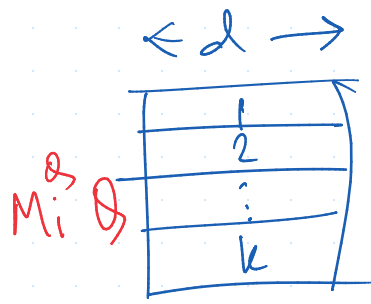
Original
query



New vector.



represents topics.



Attention head i

$$\text{output}_i = \text{Attention}(M_i^K K, M_i^V V, M_i^Q Q)$$

$$\text{Multi-Head op} = \text{Concat}(\text{op}_1, \text{op}_2, \dots, \text{op}_k) W^O$$

Multi-head attention :-

- Allows attention to focus on different aspects of the language/words.
- h -projections allow the model to attend based on different (projections) of word vectors.