

Transformer

Attention is all you need. (2017.)
(Masked) Self-Attention

~~**ELMo (2018.)**
Bi-LSTM~~

GPT(2018.)
Self-Attention (AR)

BERT(2018.)
Self-Attention (Masked token)

XLNet(2019.)
PLM (AR+AE)

Attention is all you need

Author: Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L.,
Gomez, A. N., ... & Polosukhin, I..

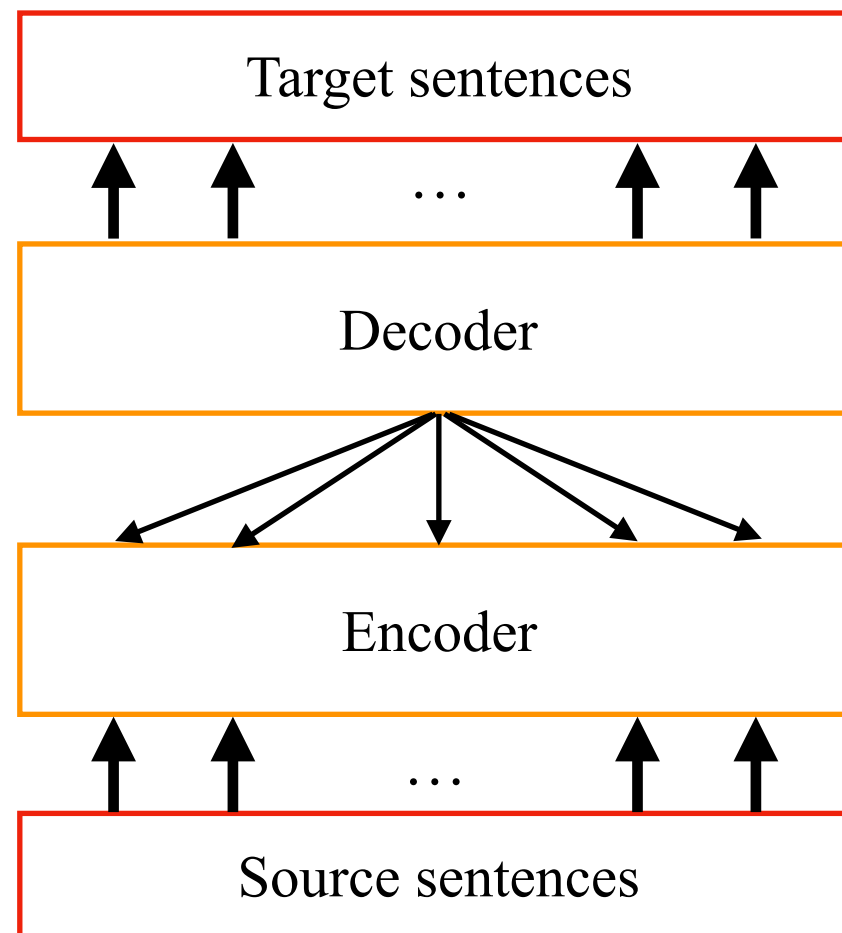
Publish: *Advances in Neural Information Processing Systems*

Pp: 5998 - 6008.

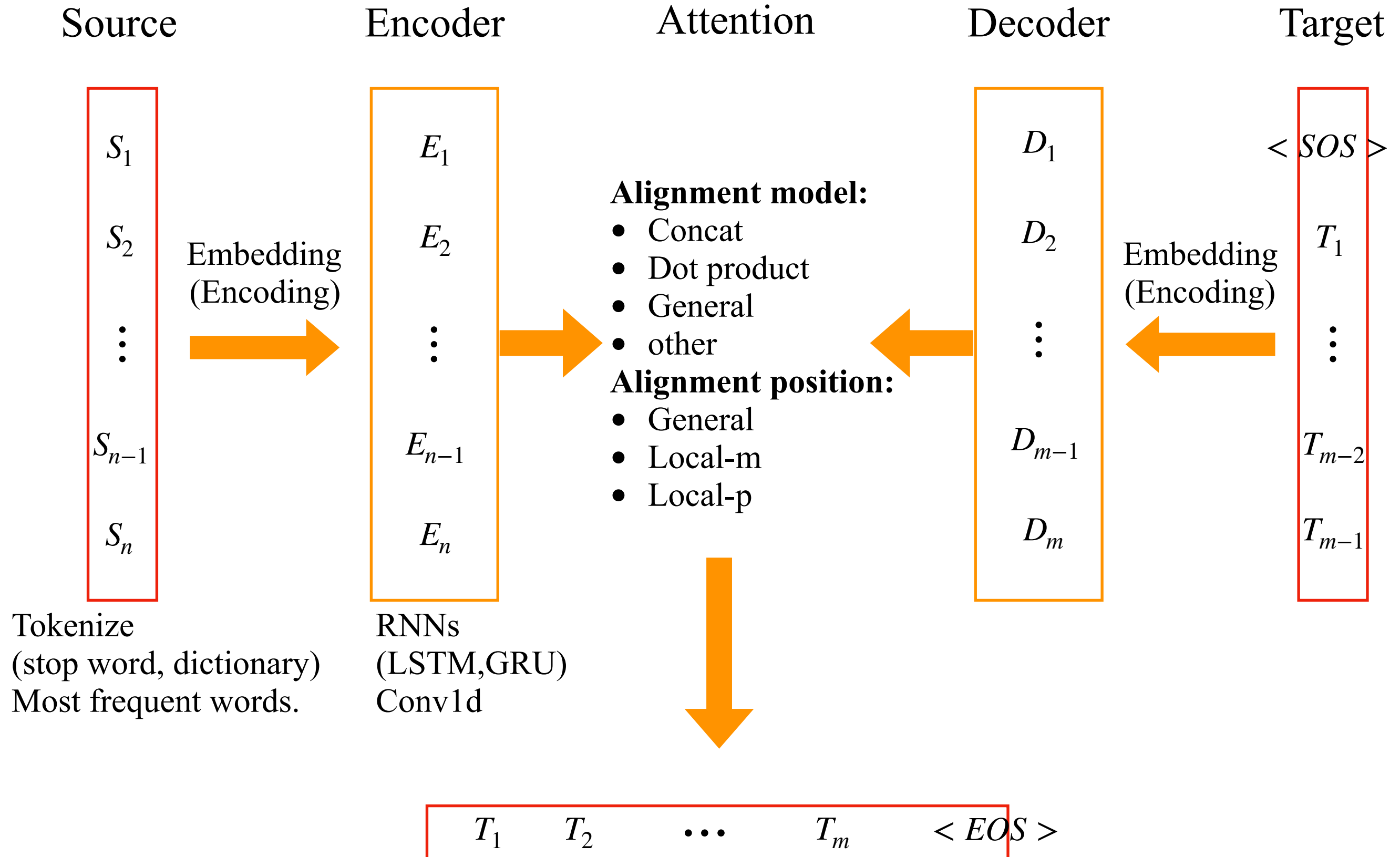
Introduction

Neural Machine Translation(NMT):

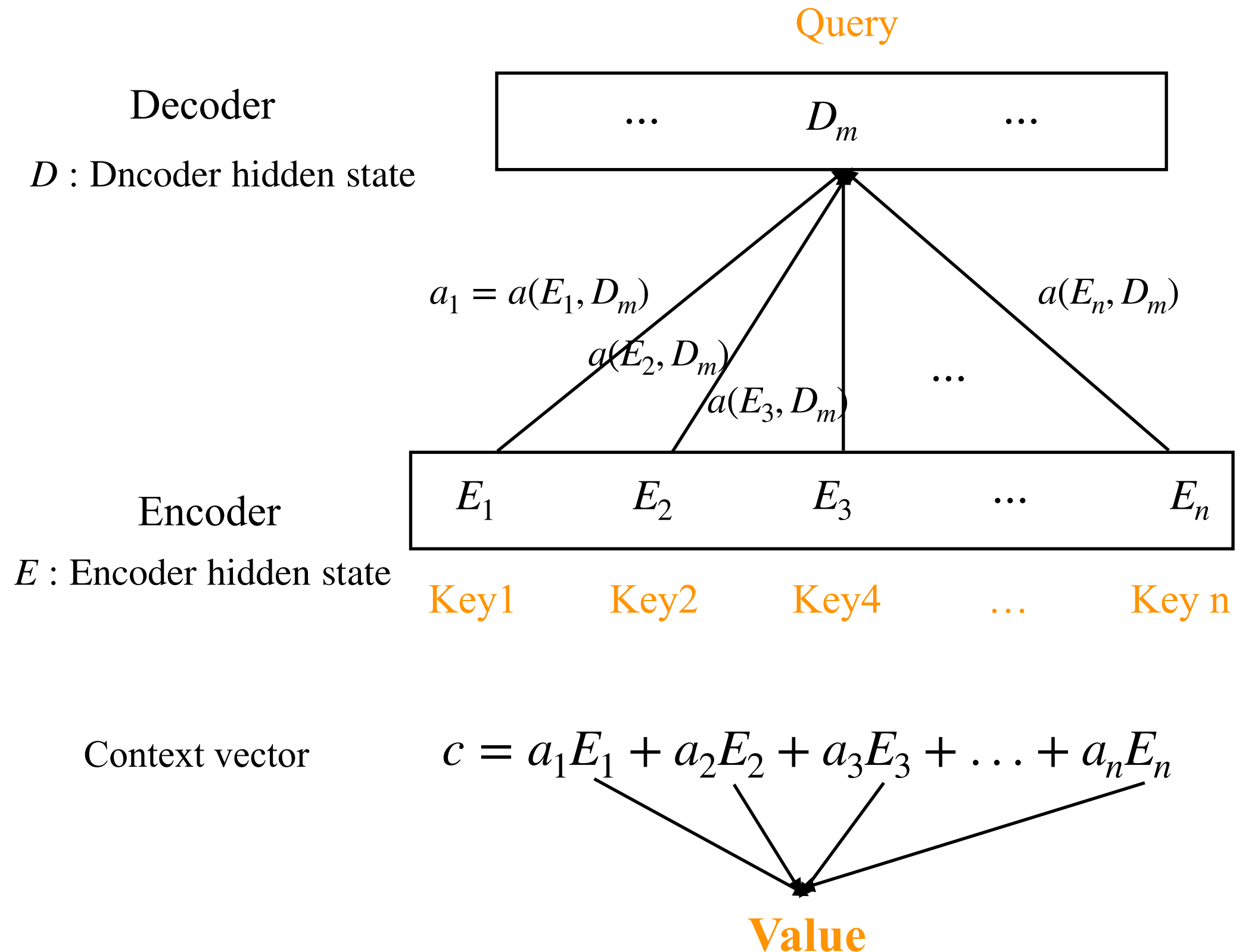
- **Statistical based:** Phrase-based + large LM (Moses)
- **NN based:** Encoder - Decoder (Seq2seq, ConvS2S, ensemble ...)



Related work



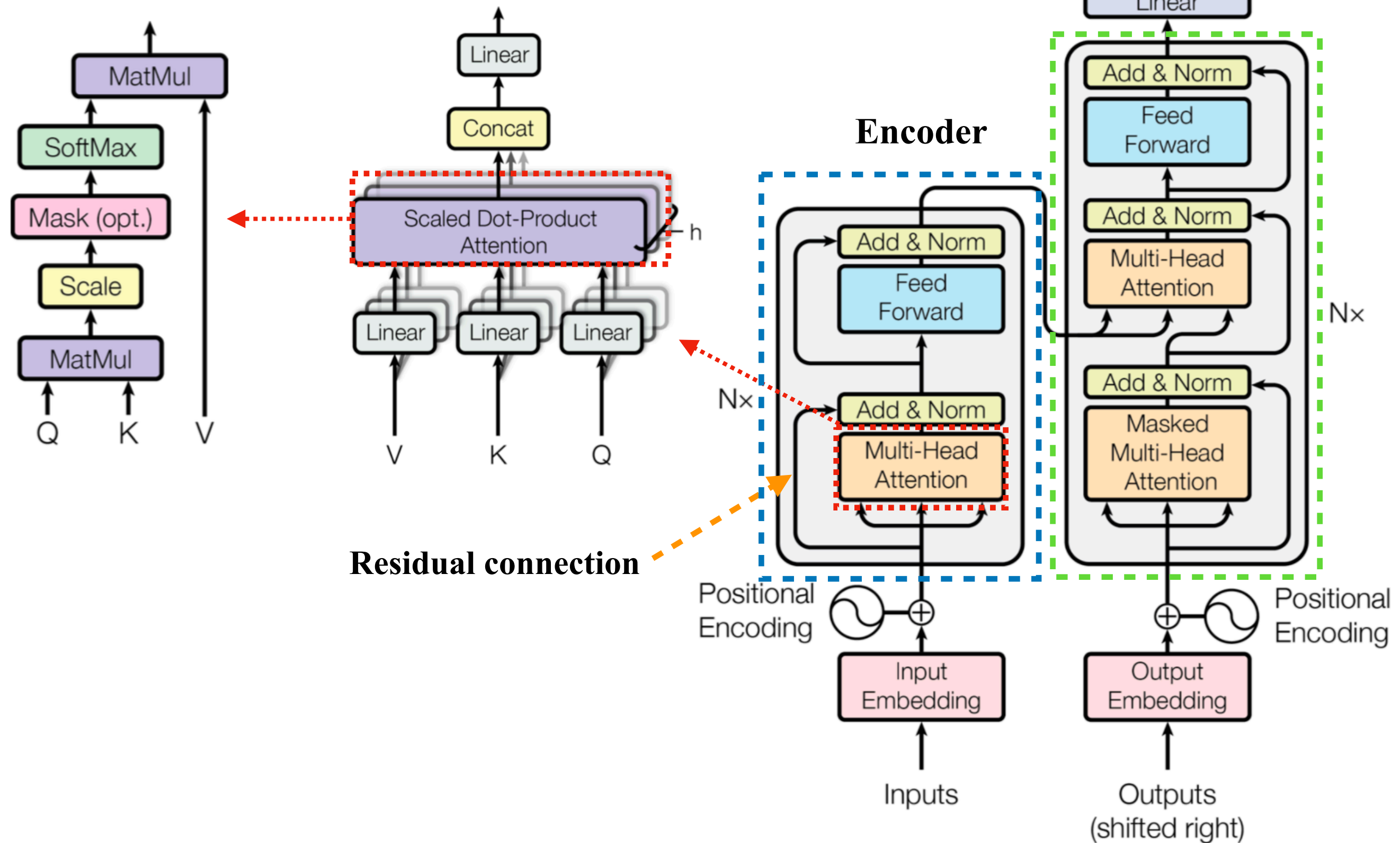
Related work



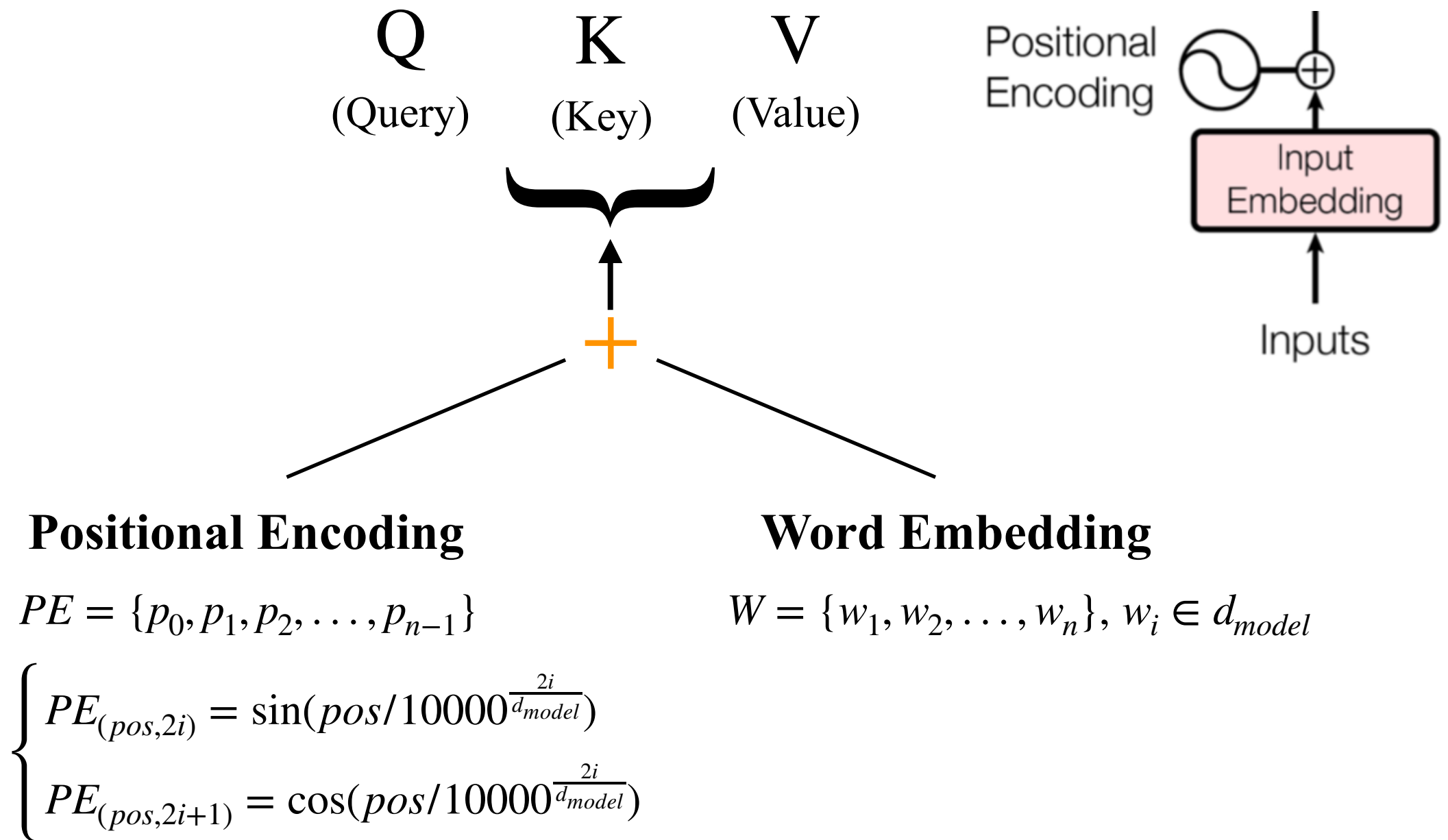
Transformer

Inner Product

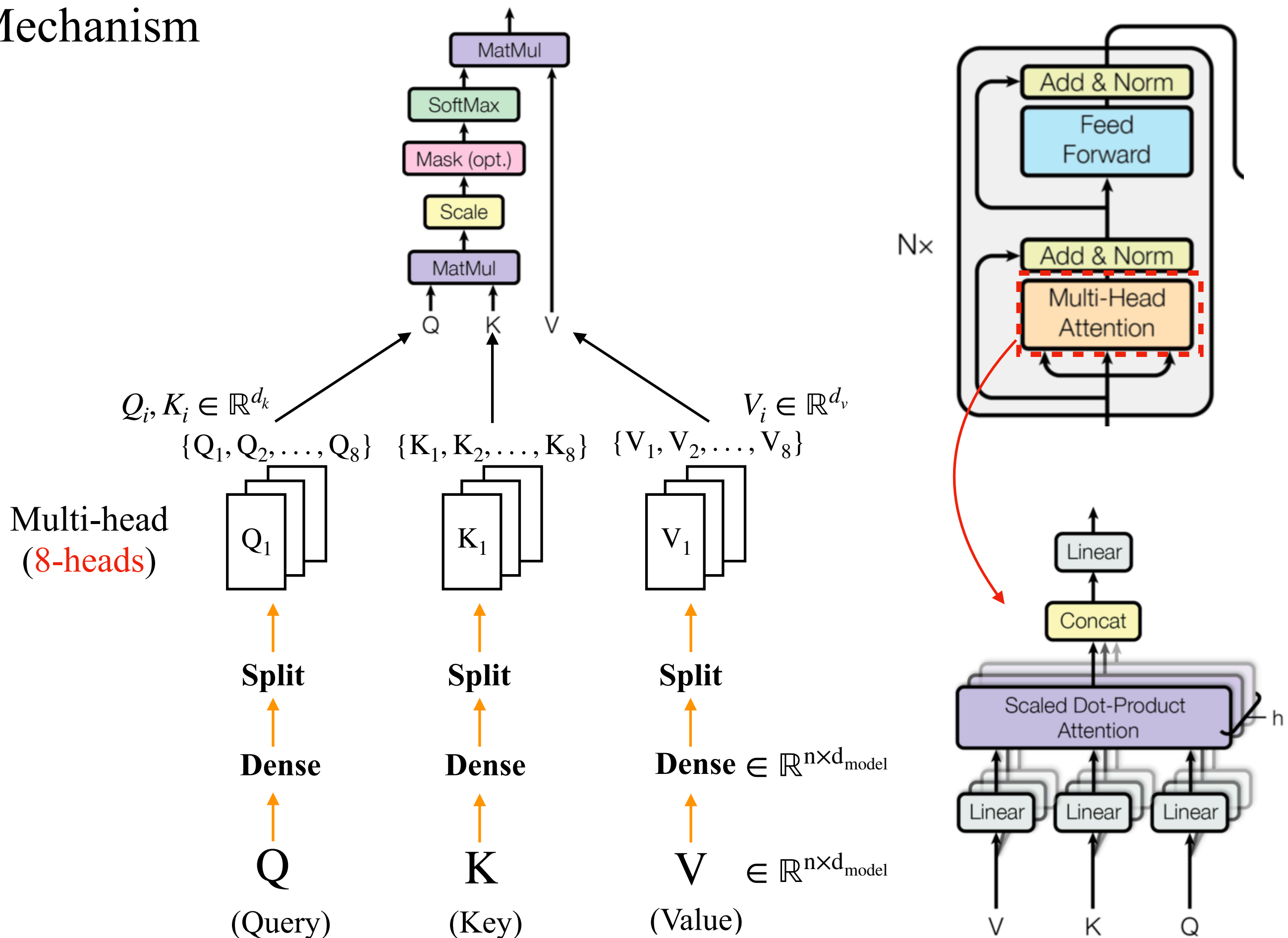
Mechanism



Mechanism



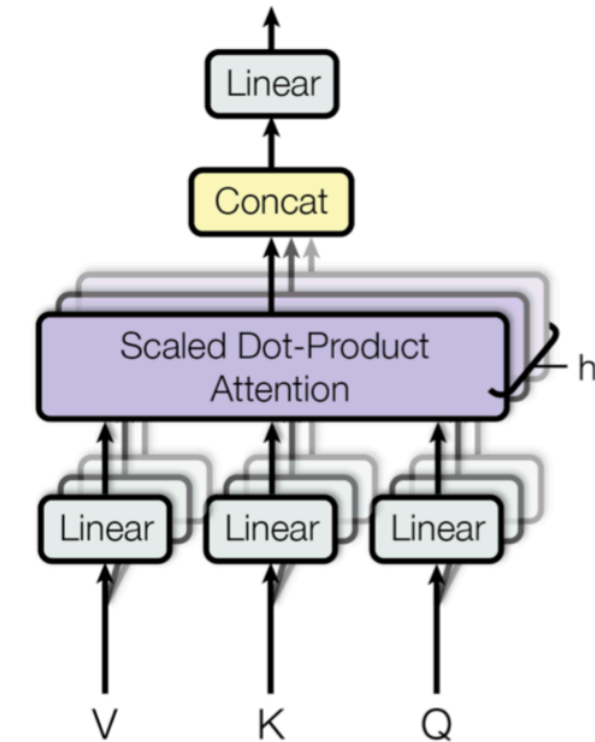
Mechanism



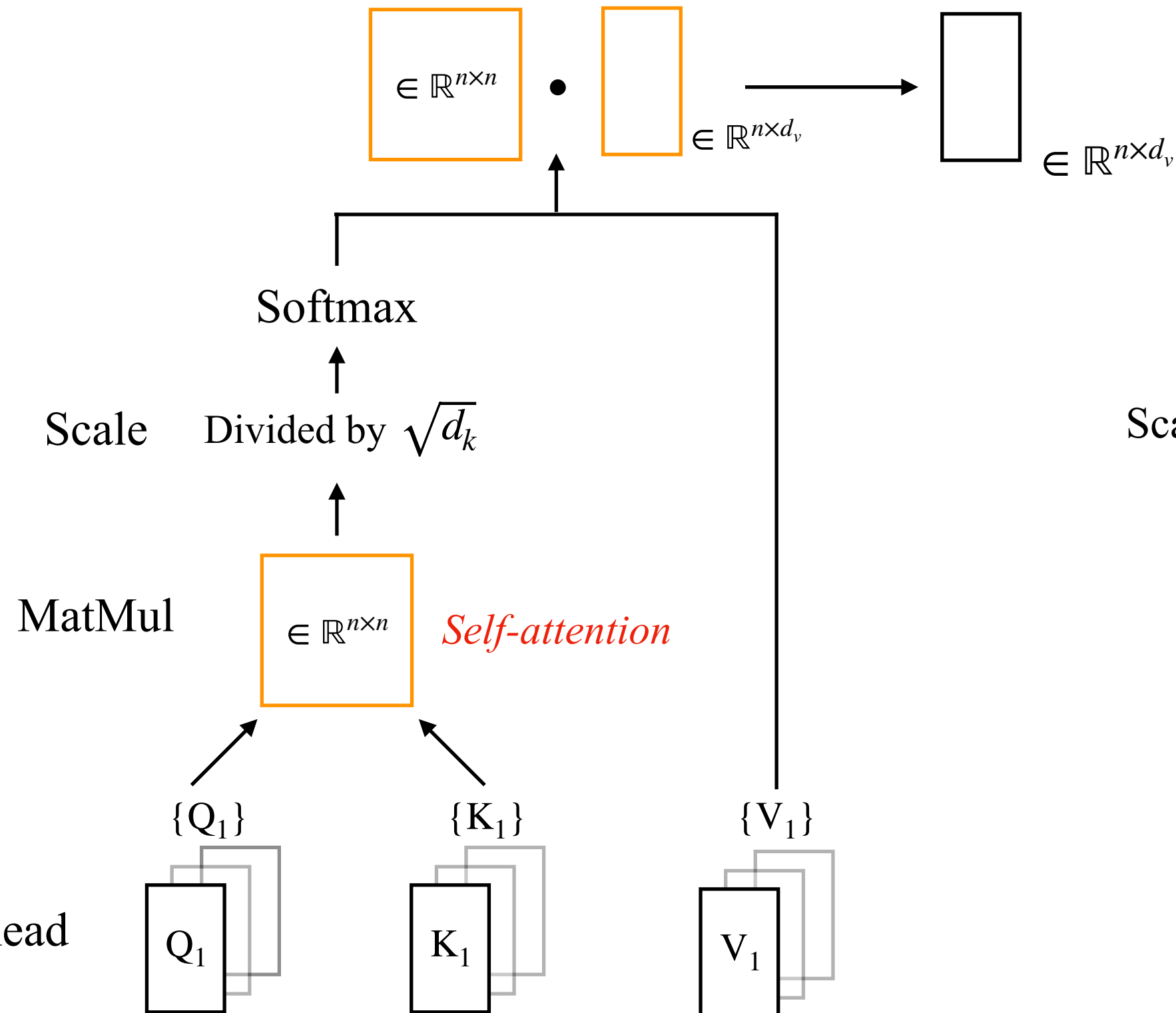
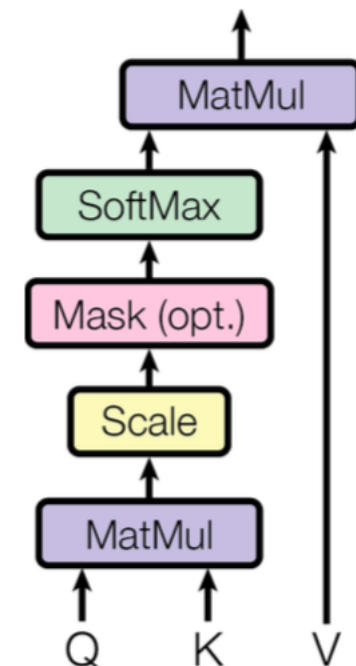
Mechanism

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Multi-heads attention

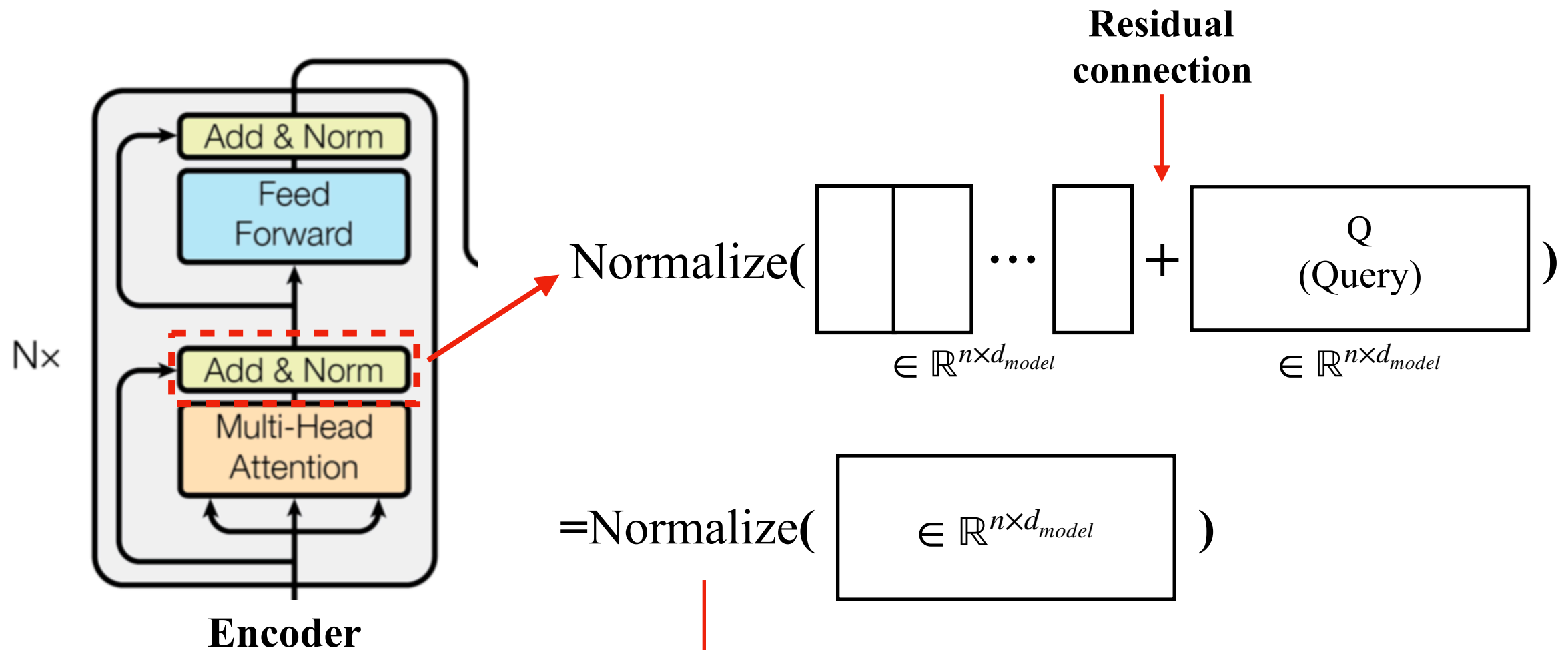


Scaled Dot-Product Attention (Self-Attention)



One-head

Mechanism



Layer Normalize: $LN(z; \alpha, \beta) = \frac{(z - \mu)}{\sigma} \odot \alpha + \beta$

Mean : $\mu^l = \frac{1}{D} \sum_{i=1}^D z_i^l$

Gains : α

Biases : β

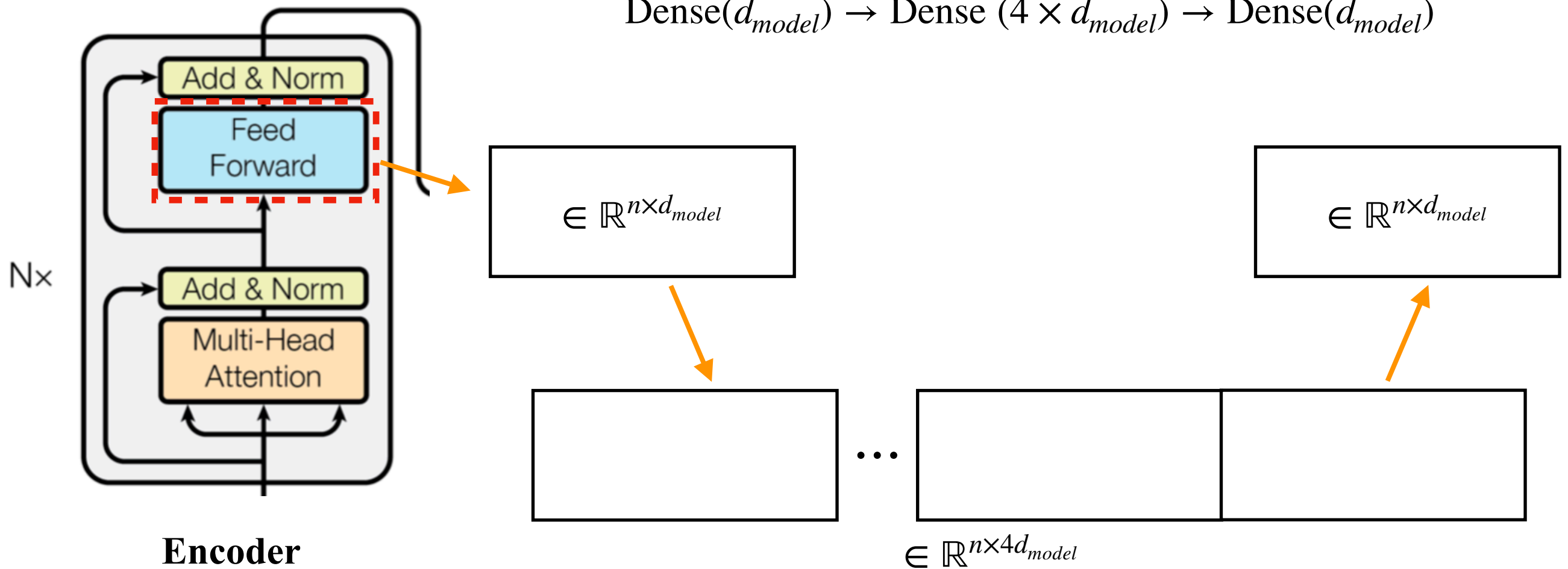
Standard Deviation : $\sigma^l = \sqrt{\frac{1}{D} \sum_{i=1}^D (z_i^l - \mu^l)^2}$

Mechanism

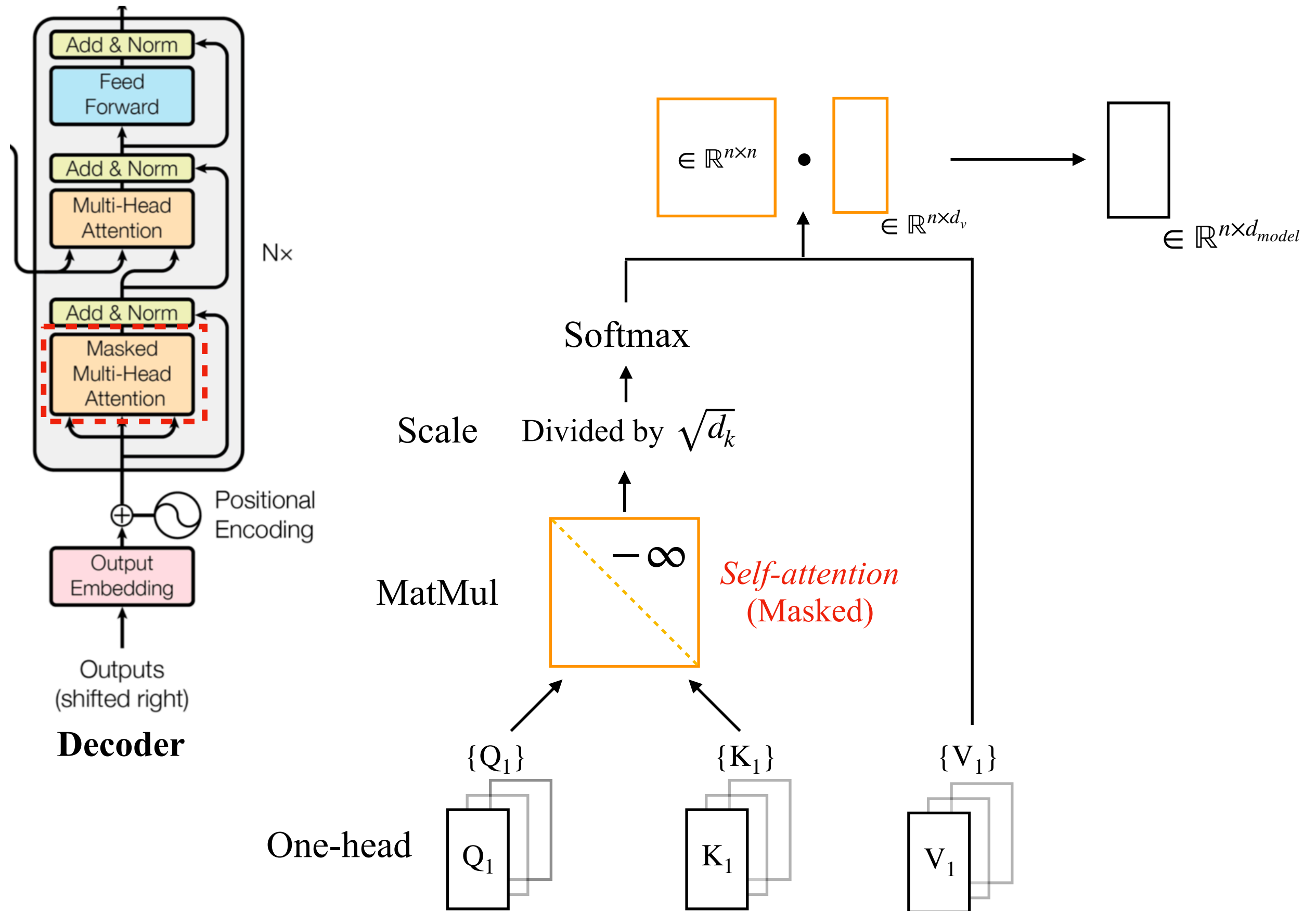
Feed Forward (Dense 、Conv1d)

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

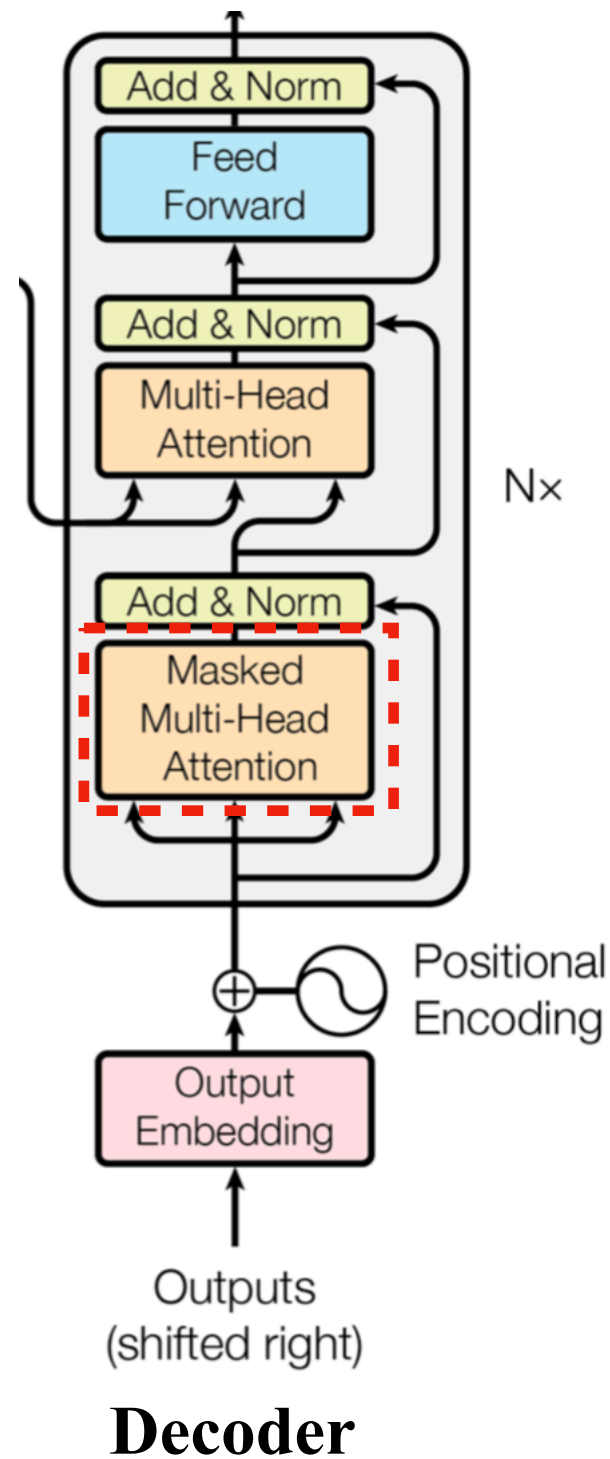
$$\text{Dense}(d_{\text{model}}) \rightarrow \text{Dense}(4 \times d_{\text{model}}) \rightarrow \text{Dense}(d_{\text{model}})$$



Mechanism



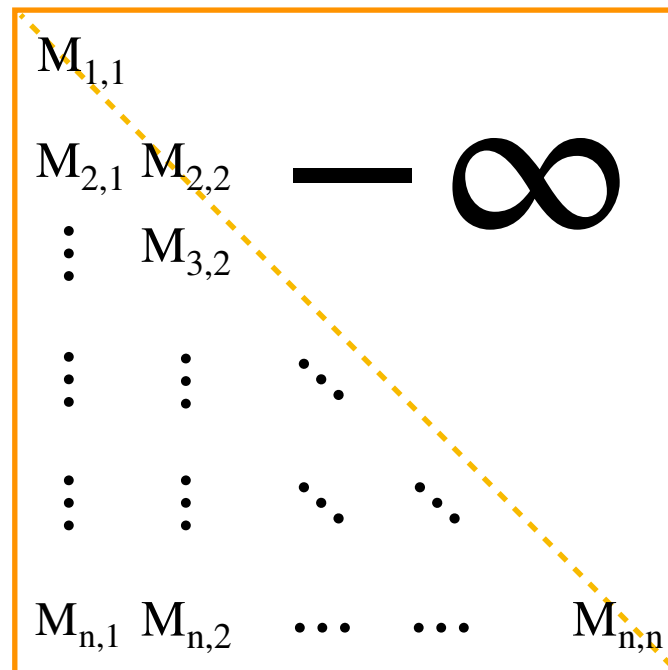
Mechanism



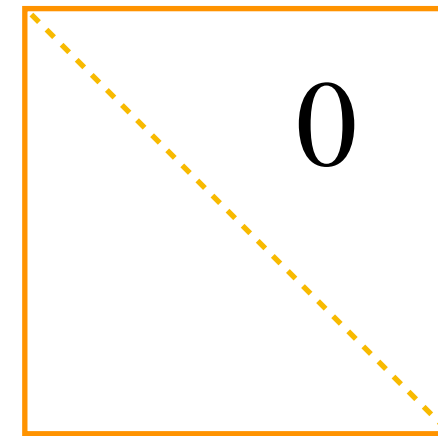
$N \times$ We need to prevent leftward information flow in the decoder
 $M_{1,1}$: The similarity between Q_1 and K_1

One-head

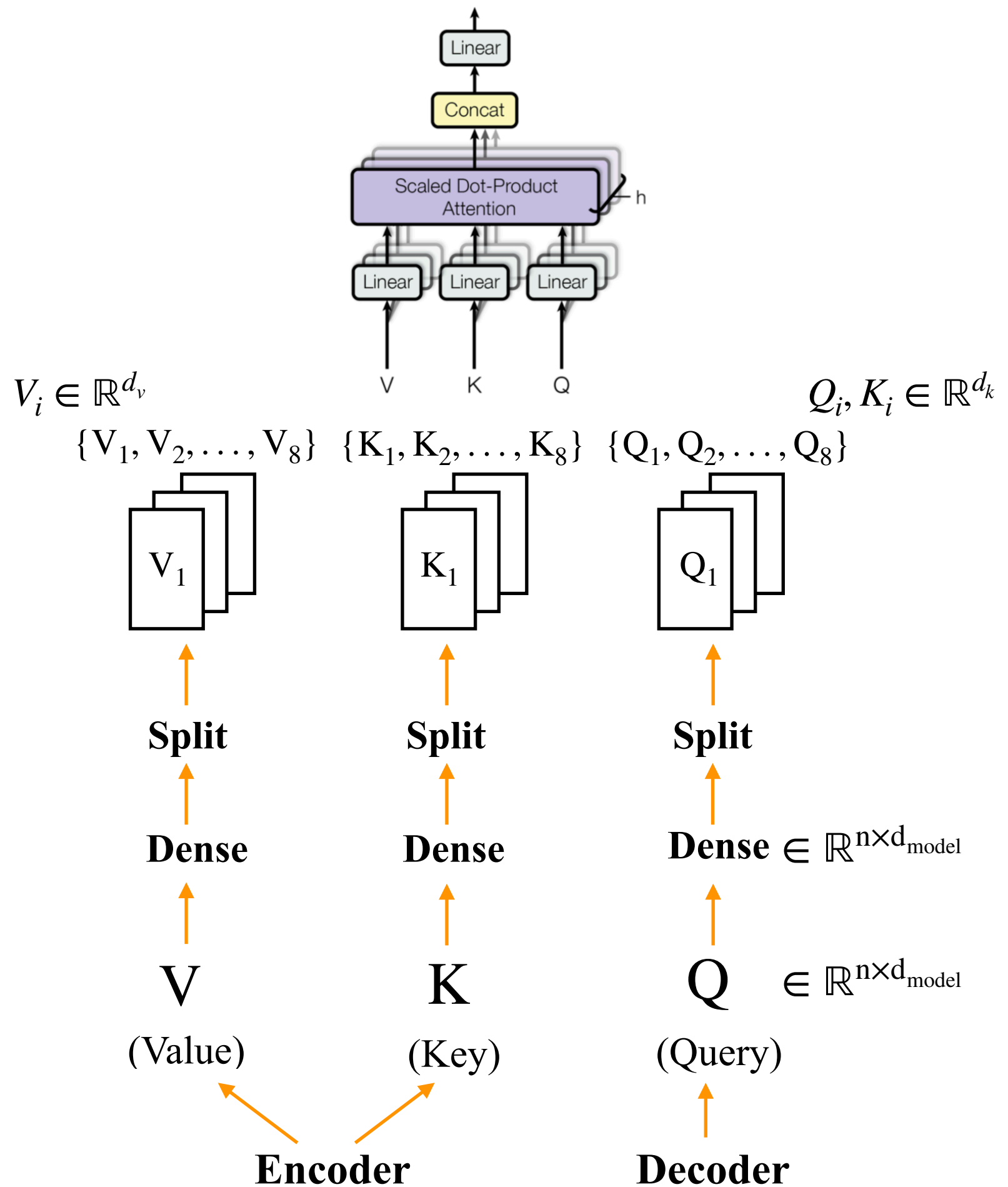
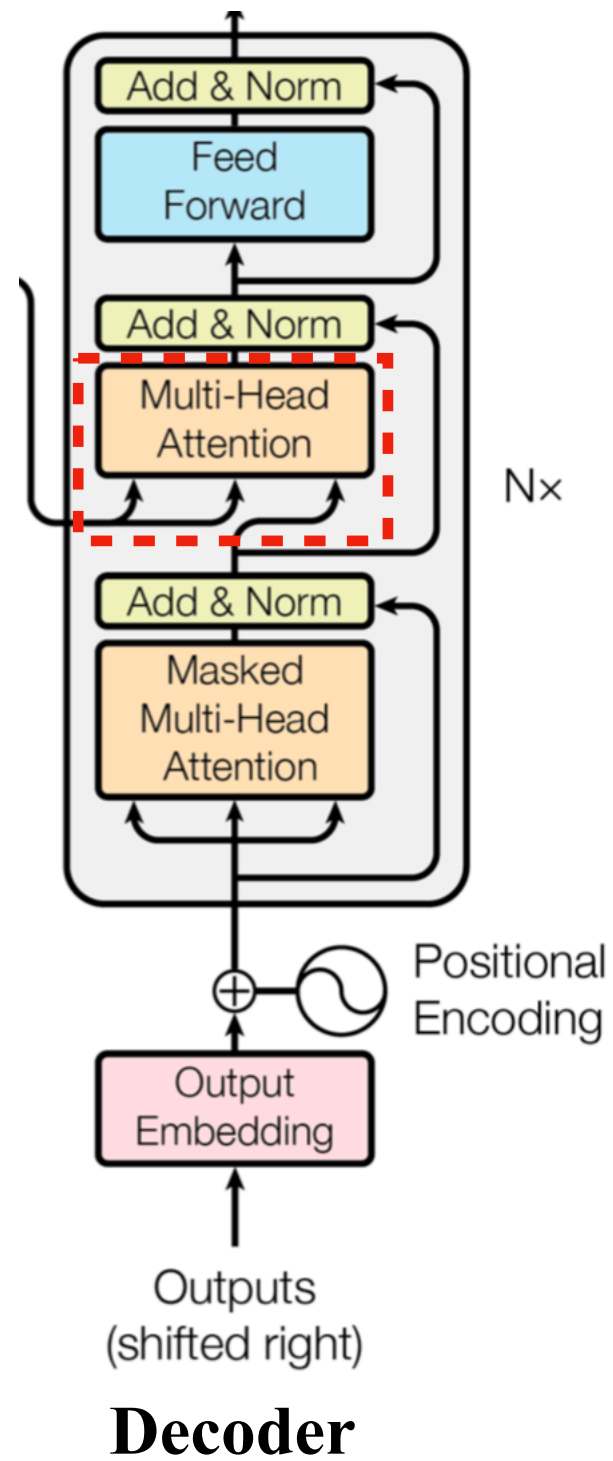
Scale Divided by $\sqrt{d_k}$ Softmax



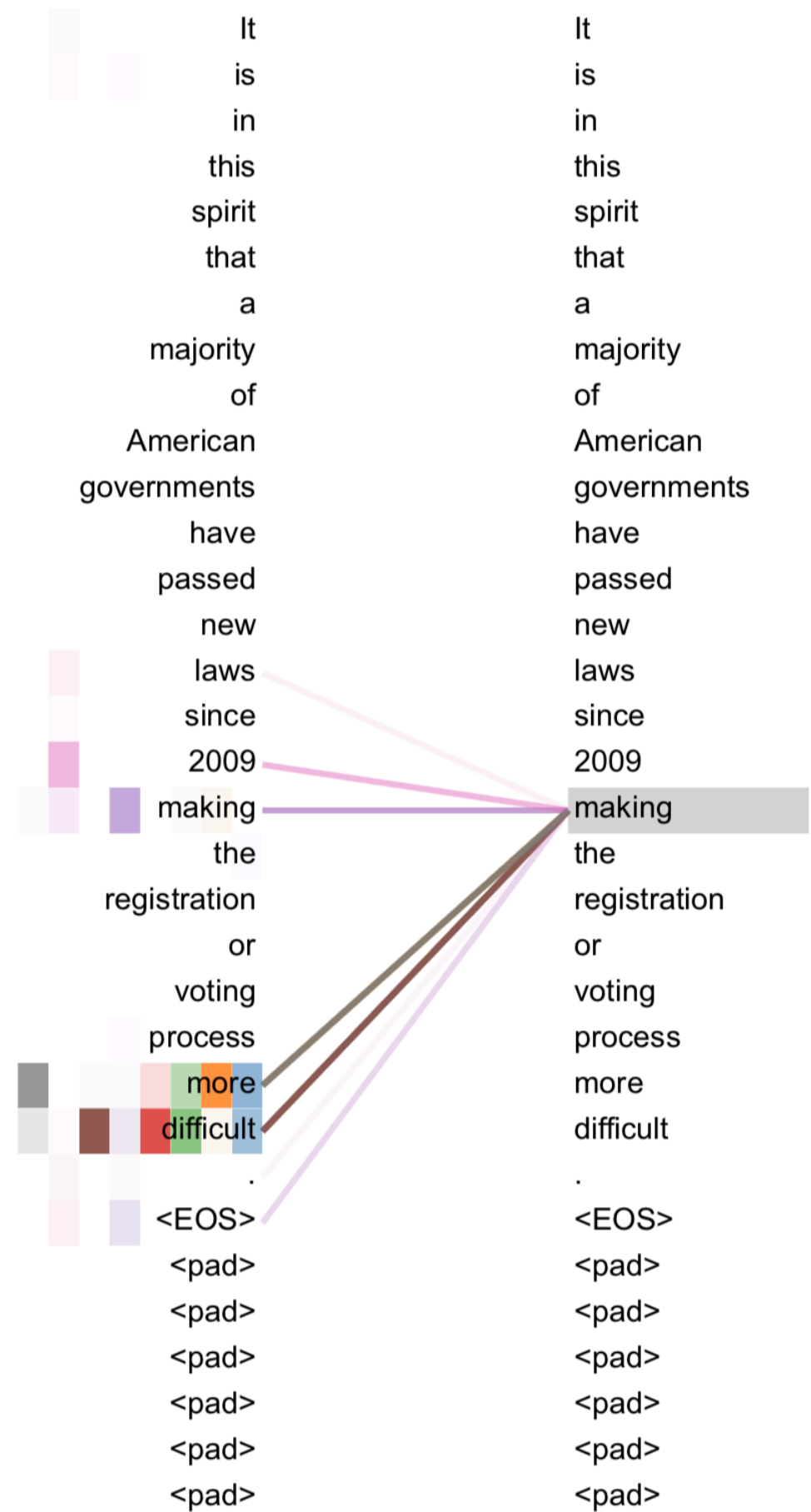
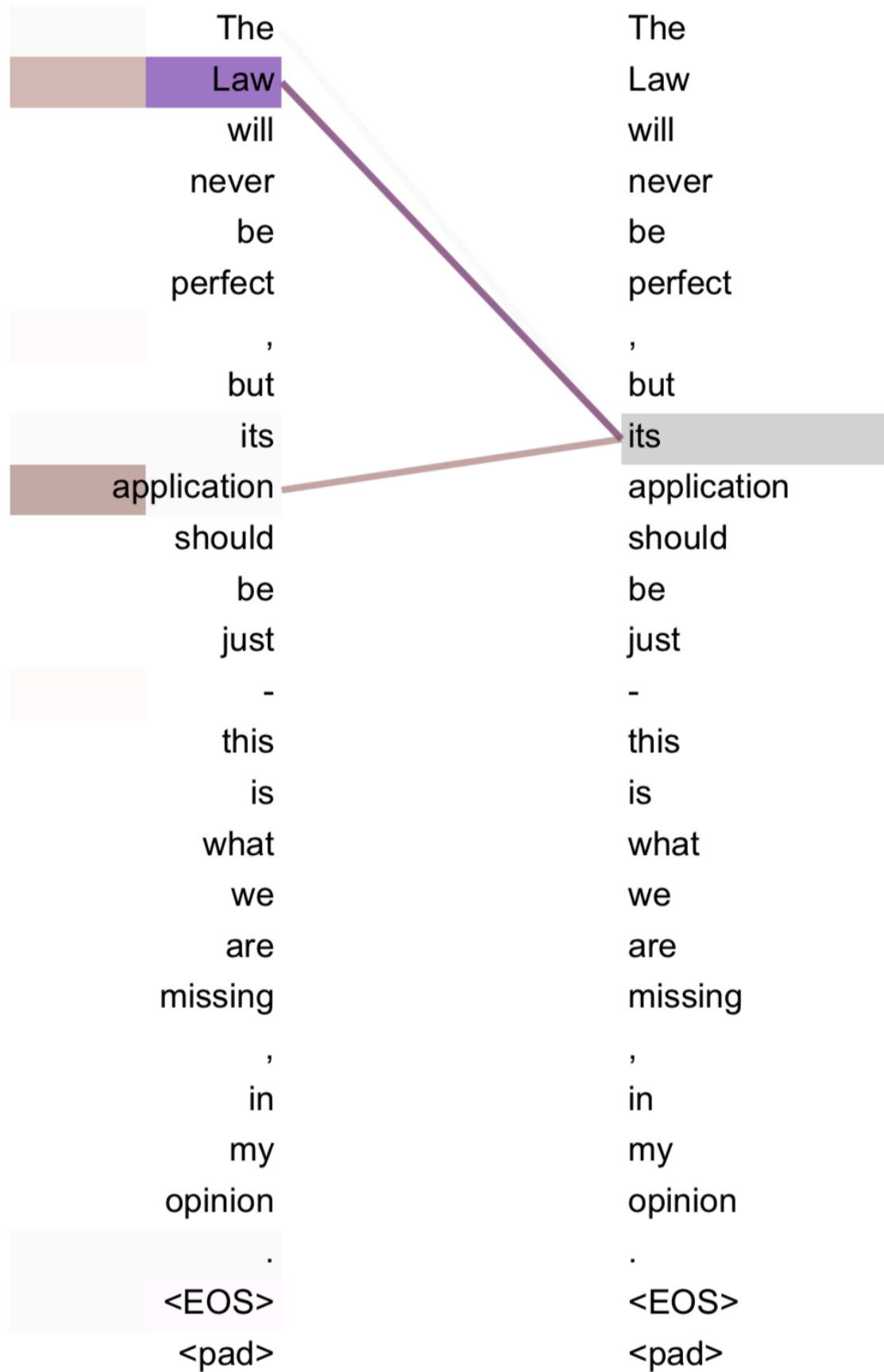
*Self-attention
(Masked)*



Mechanism



Evaluation



- **Attention mechanism:**

<https://reurl.cc/NzNne>

- **Attention is all you need:**

<https://reurl.cc/nWlX1>

- **BERT:**

<https://reurl.cc/76bld>

- **XLNet:**

<https://reurl.cc/50baV>