

Attention

is all you need

Ascii Paper Reading Group

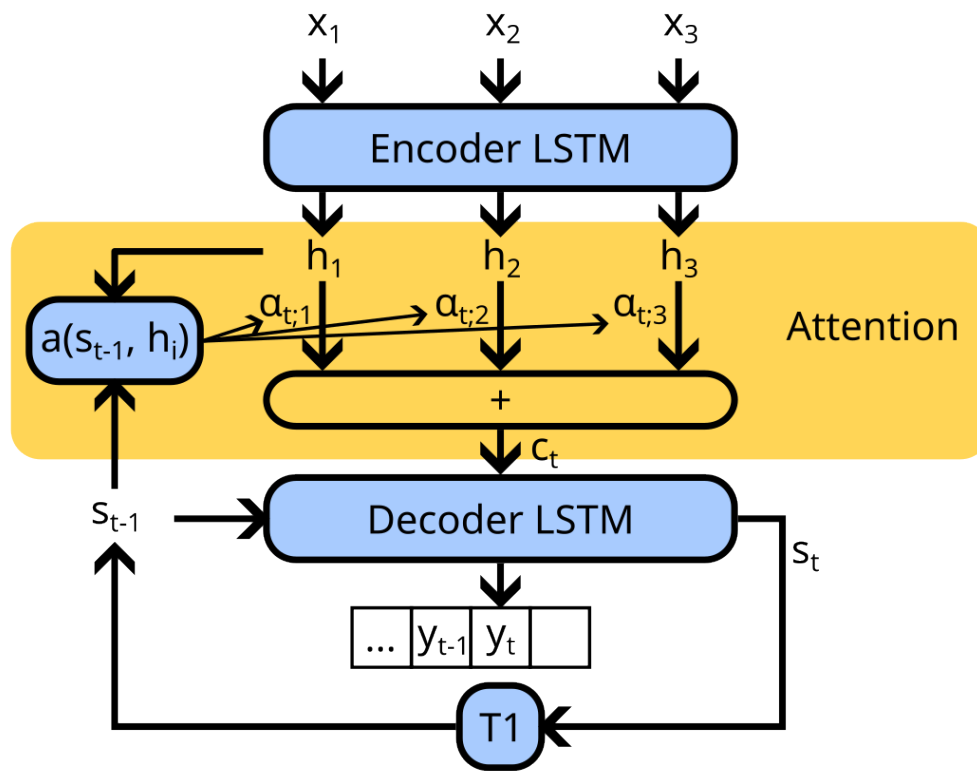


Goal



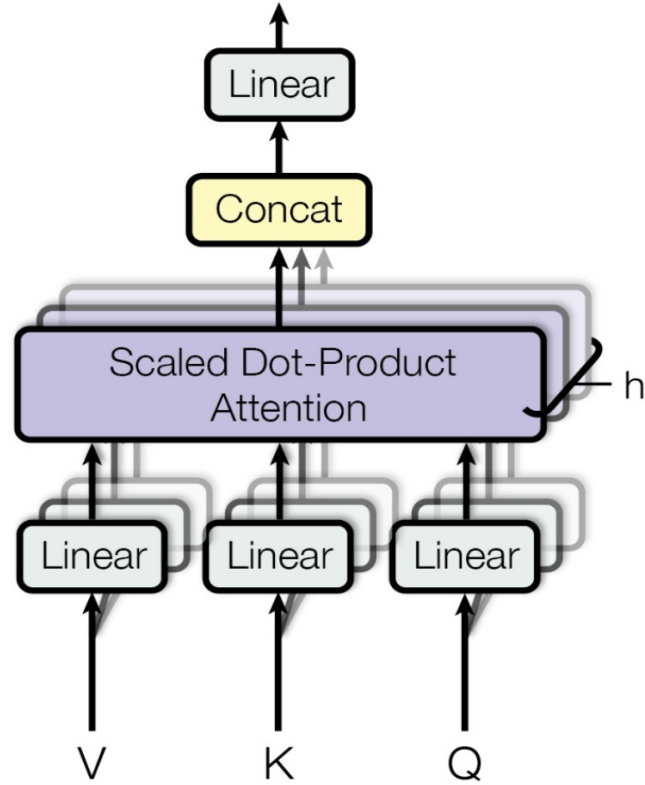
- **Machine Translation**
 - requires alignment of source with output sentence
- **Parallel calculation**

Previous Works



Alignment using

Multi Head Attention



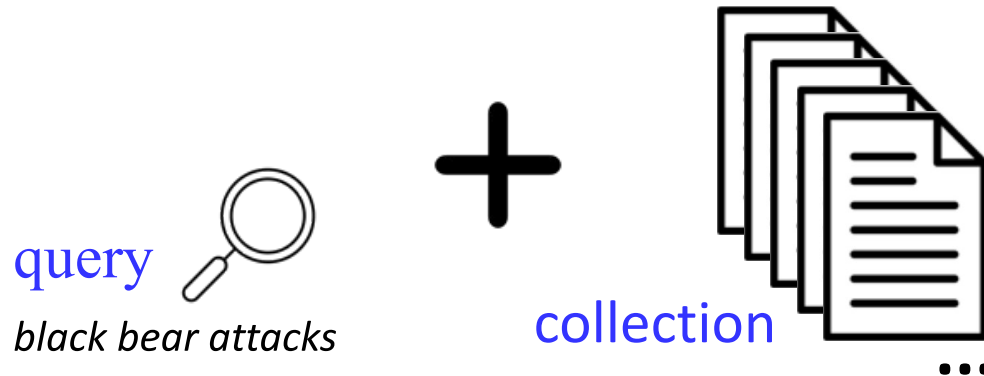
Focus: Ad hoc Retrieval

Given: query q

collection of texts

Return: a ranked list of k texts $d_1 \dots d_k$

Maximizing: a metric of interest

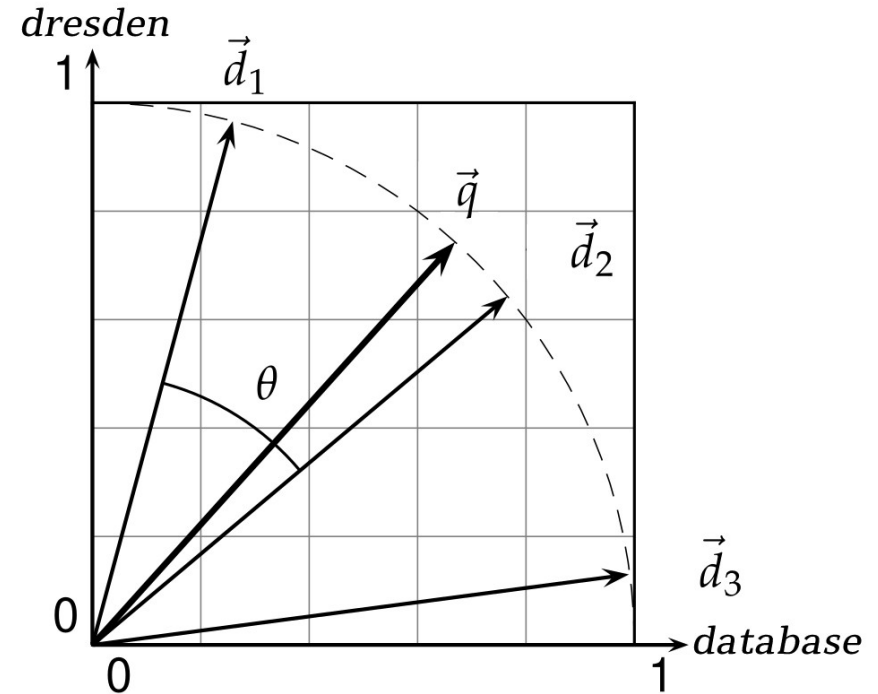


metric: 0.66

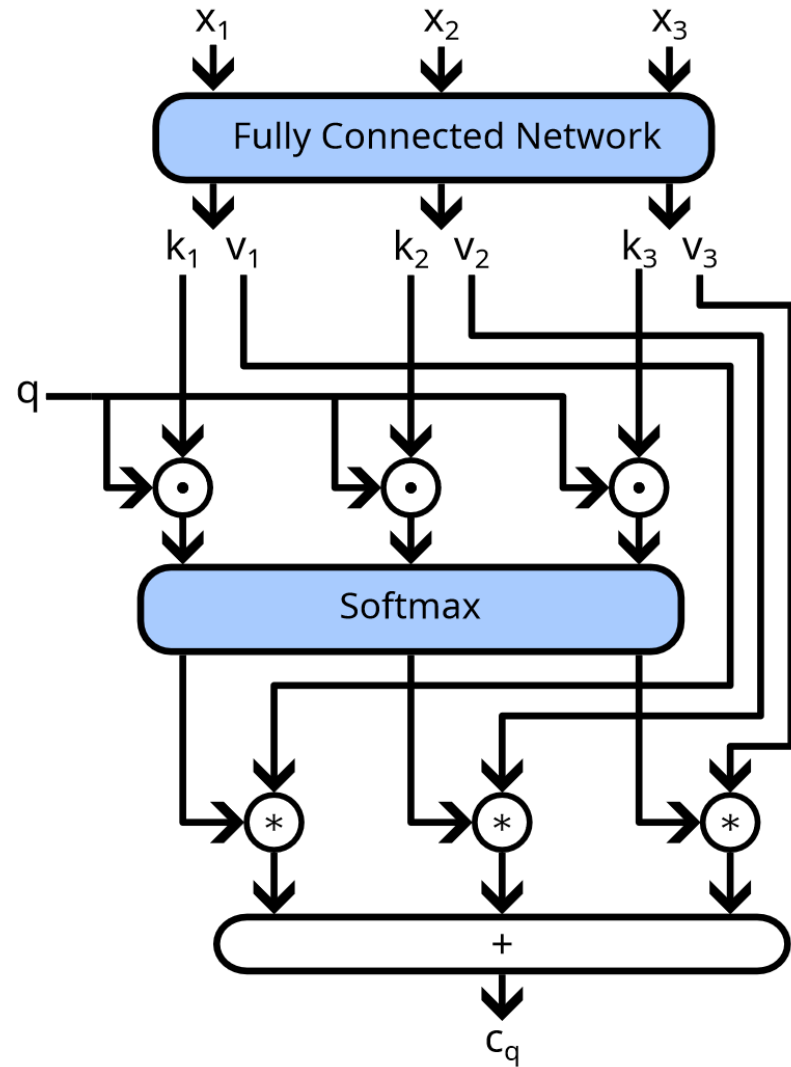
Vector Space Model

- queries q + documents d
- represented as vectors
- use cosine similarity between query and documents for ranking

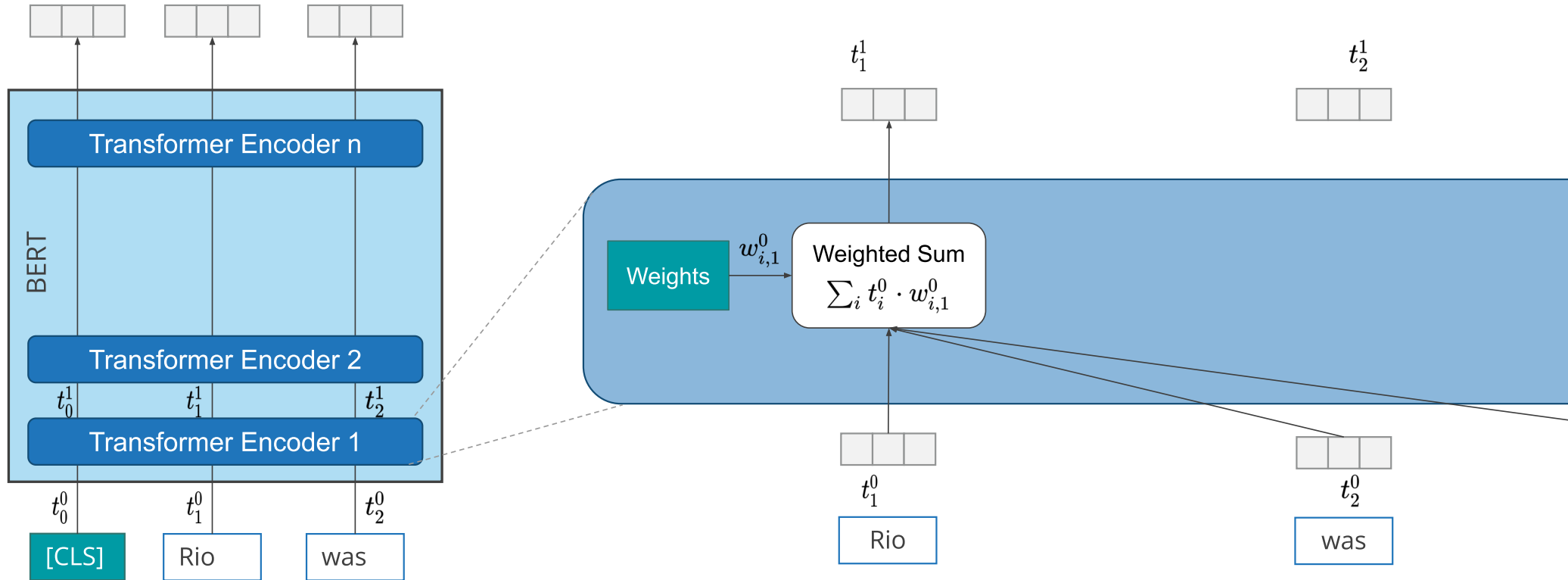
$$\text{sim}(d_1, d_2) = \cos \theta = \frac{\vec{d}_1 \cdot \vec{d}_2}{|\vec{d}_1| |\vec{d}_2|}$$



Dot Product Attention

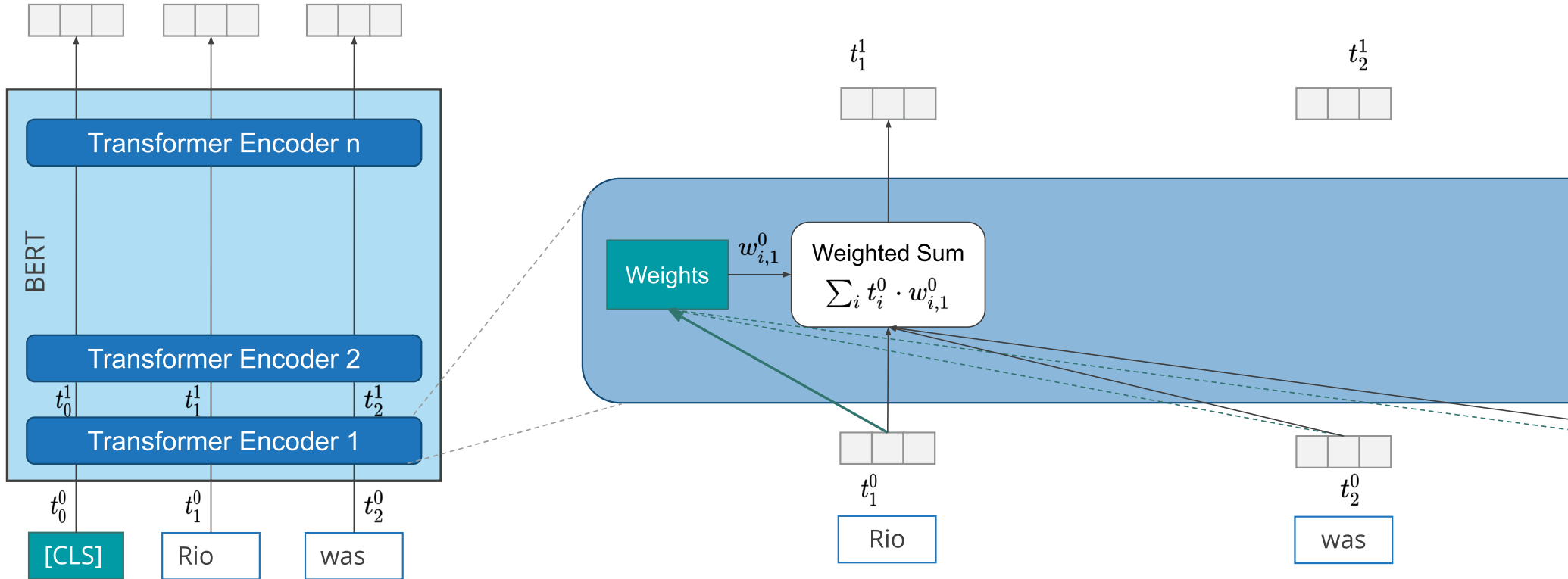


Inside BERT - Attention



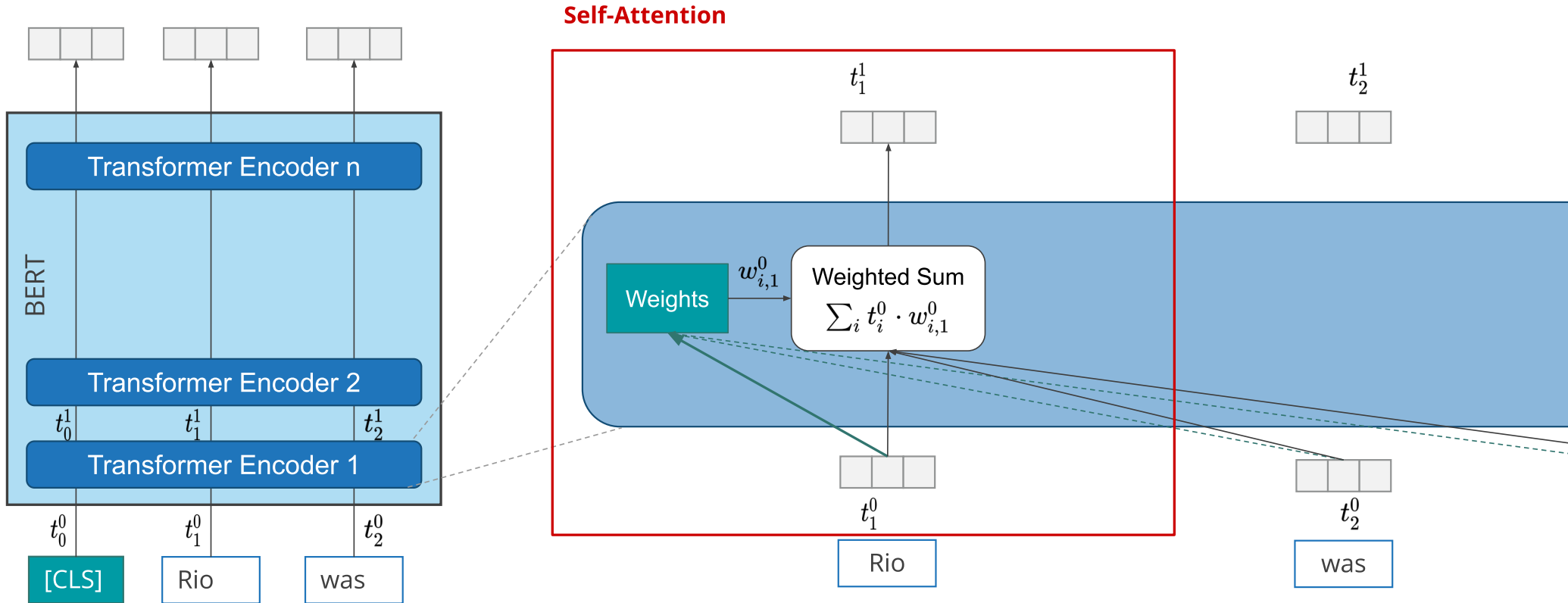
- Weights determined by current vector (mini neural network is trained here)
- Context of all other words influences the output vector of token t_i

Inside BERT - Attention



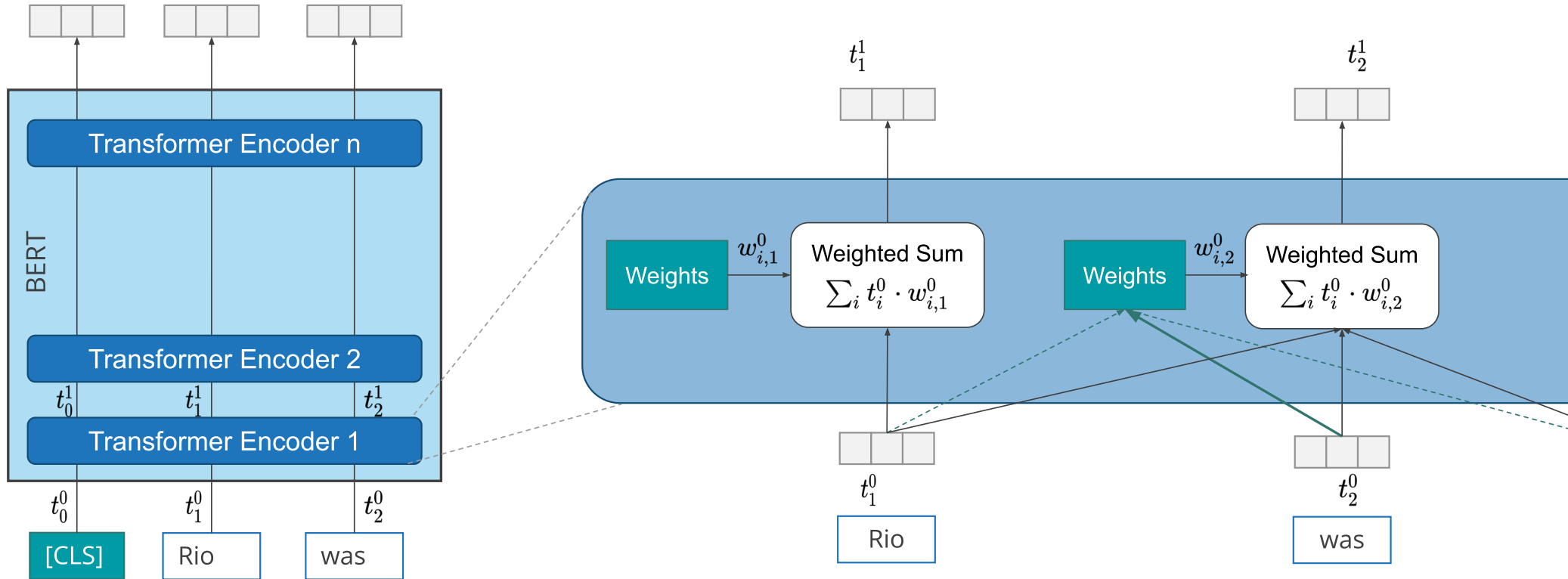
- Weights determined by current vector (mini neural network is trained here)
- Context of all other words influences the output vector of token t_i

Inside BERT - Attention



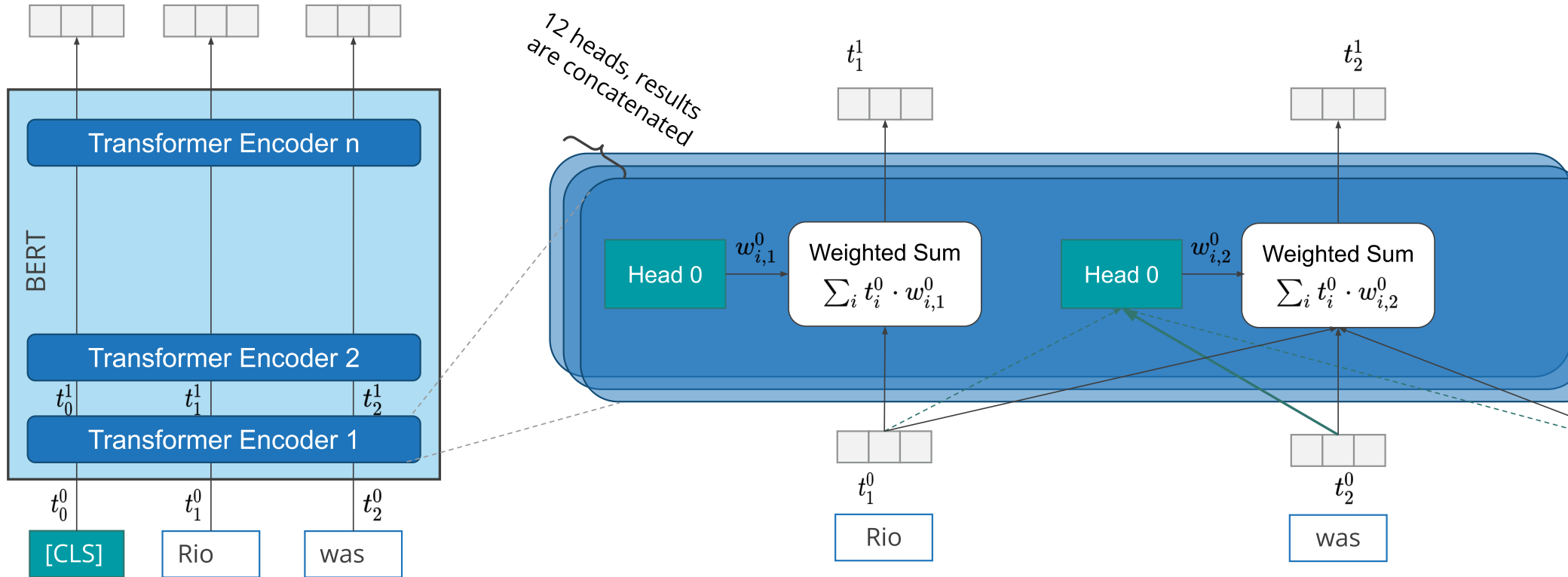
- Weights determined by current vector (mini neural network is trained here)
- Context of all other words influences the output vector of token t_i

Inside BERT - Attention



- Weights determined by current vector (mini neural network is trained here)
- Context of all other words influences the output vector of token t_i

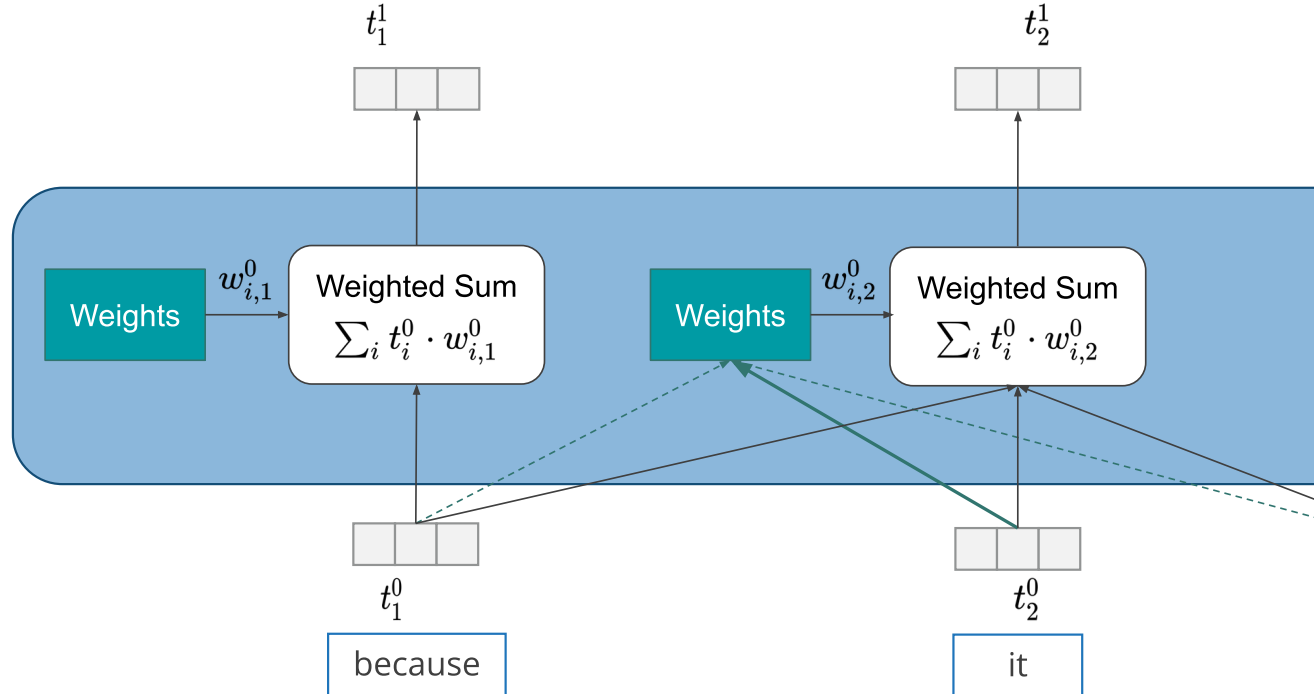
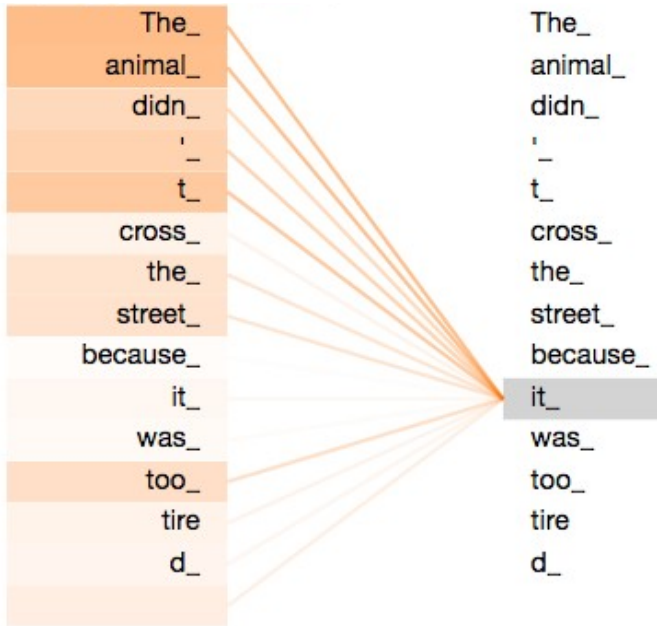
Inside BERT - Attention



- Weights determined by current vector (mini neural network is trained here)
- Context of all other words influences the output vector of token t_i

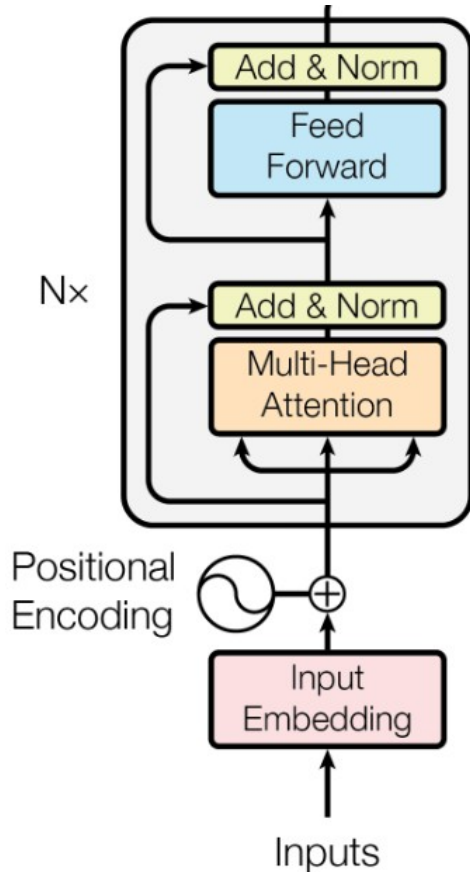
Inside BERT - Attention

Visualization of weights of one attention head:



- Weights determined by current vector (mini neural network is trained here)
- Context of all other words influences the output vector of token t_i

Encoder



- Repeated self attention
- Skip connections for gradient flow
- Normalization for
 - gradient flow
 - Cosine distance

Positional Encoding

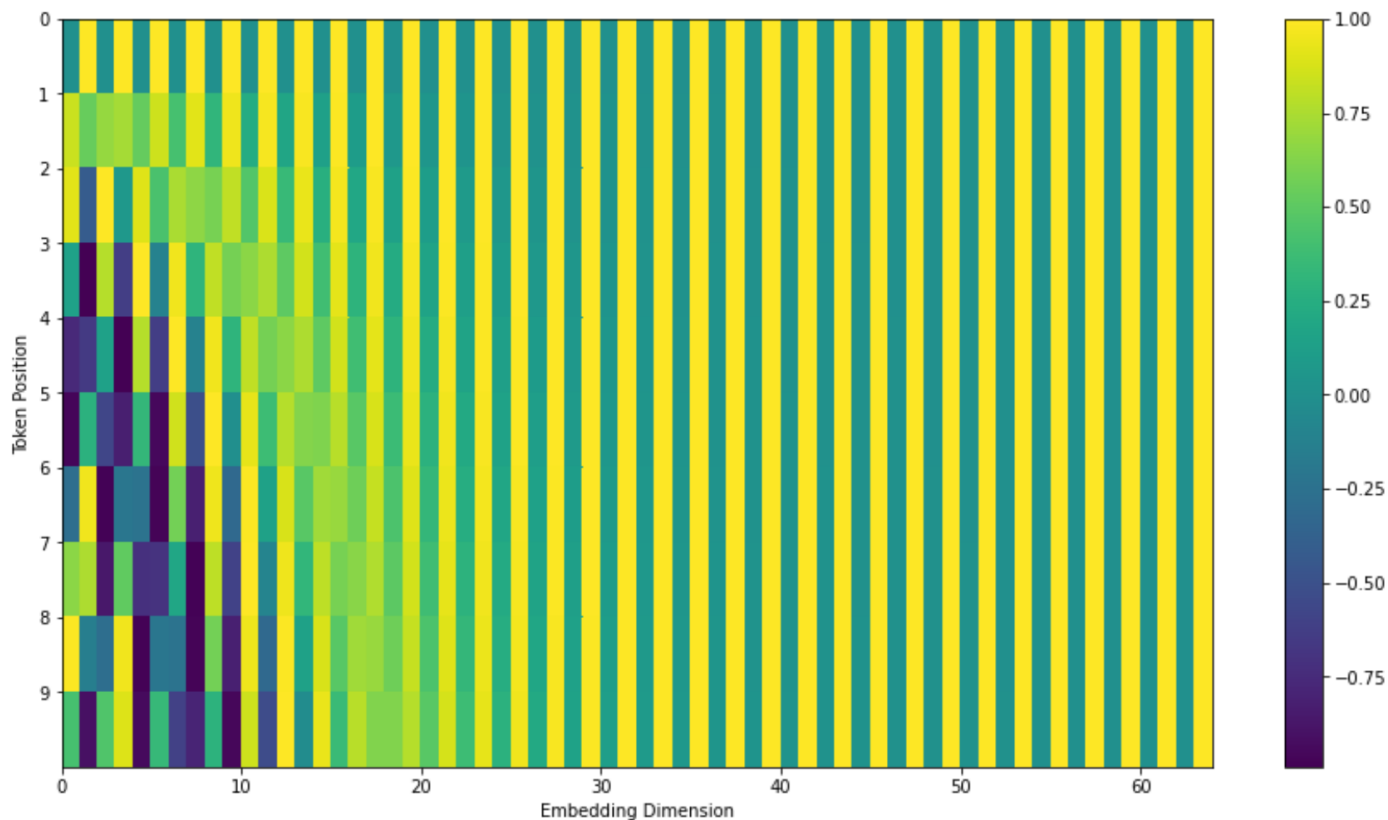
- Problem: Dot-Product attention doesn't consider distance of words in a document.
- Add a positional encoding to each token

$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

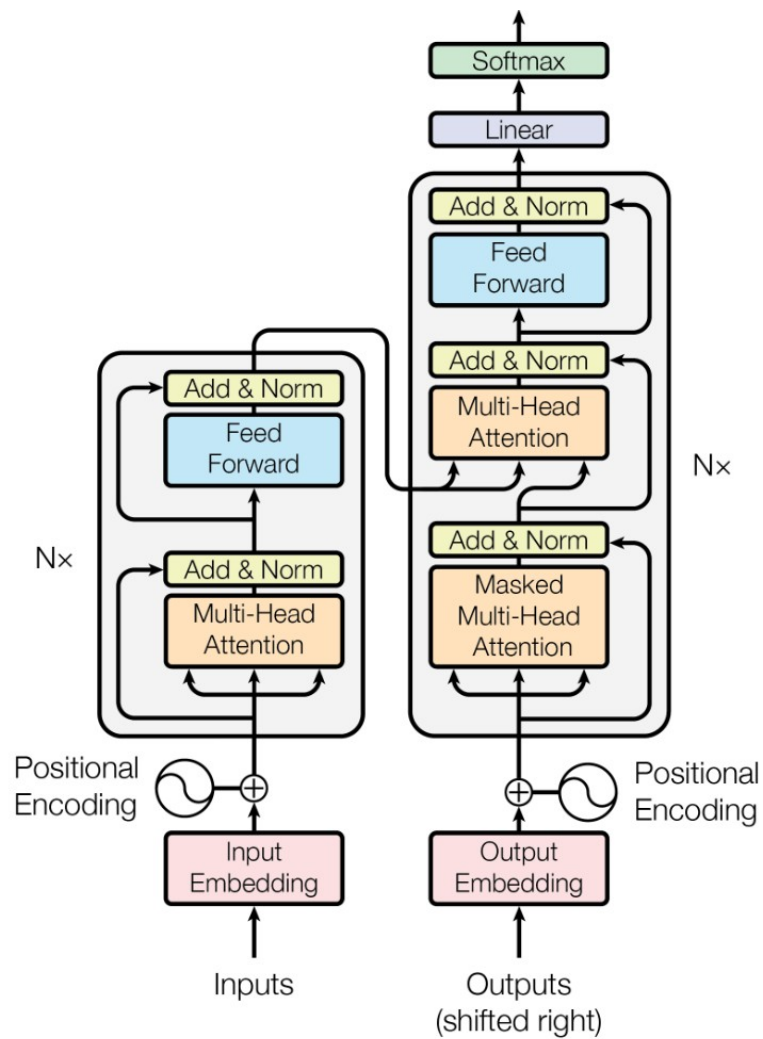
$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

i ... dimension

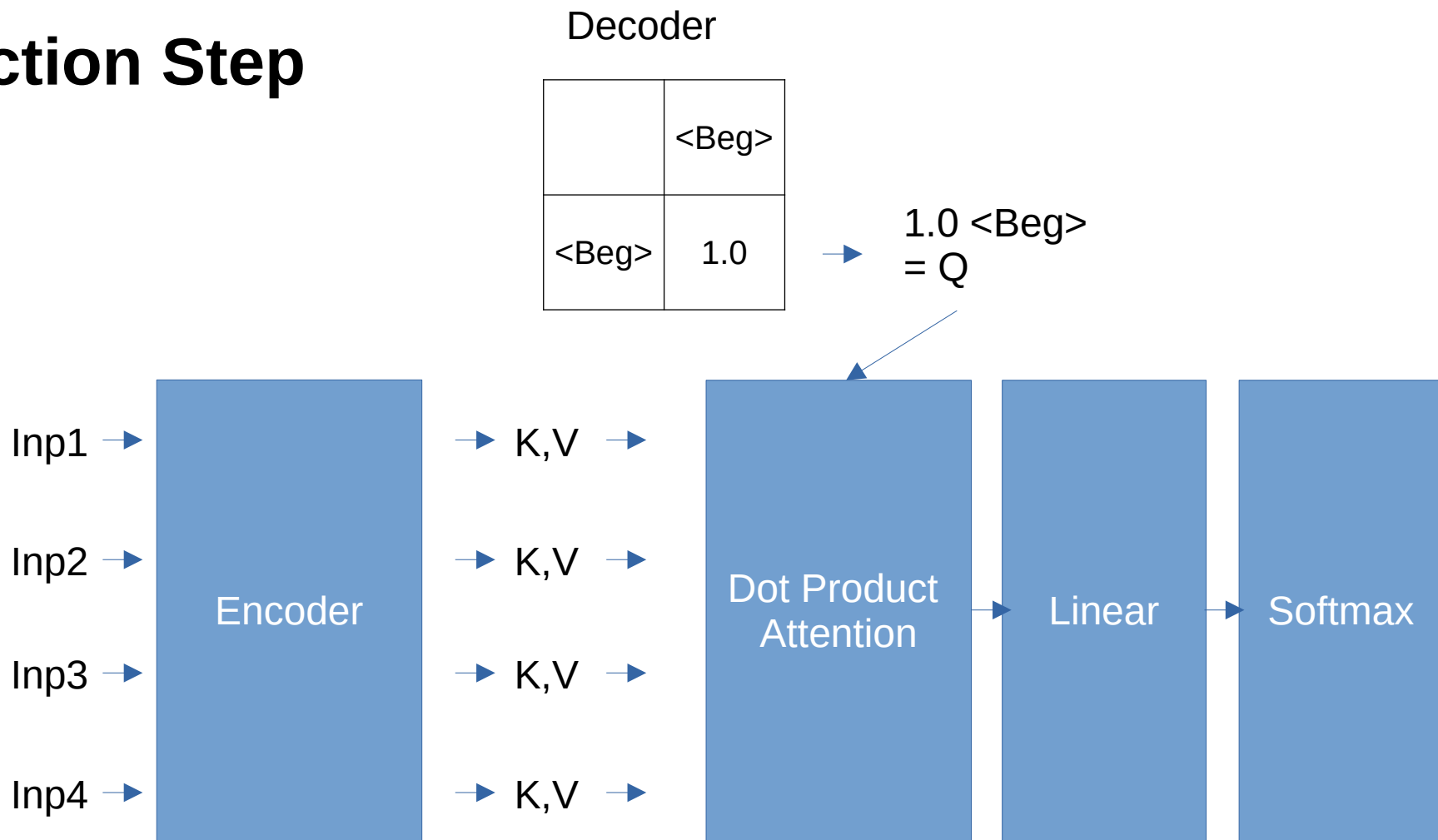
Positional Encoding



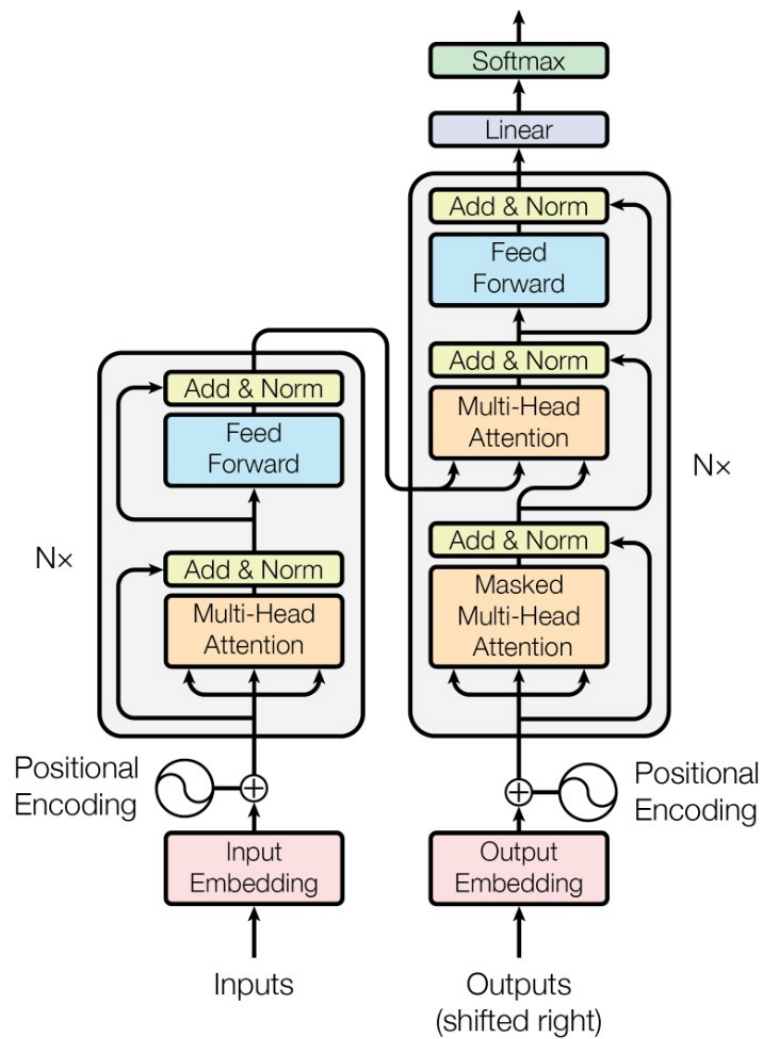
Decoder



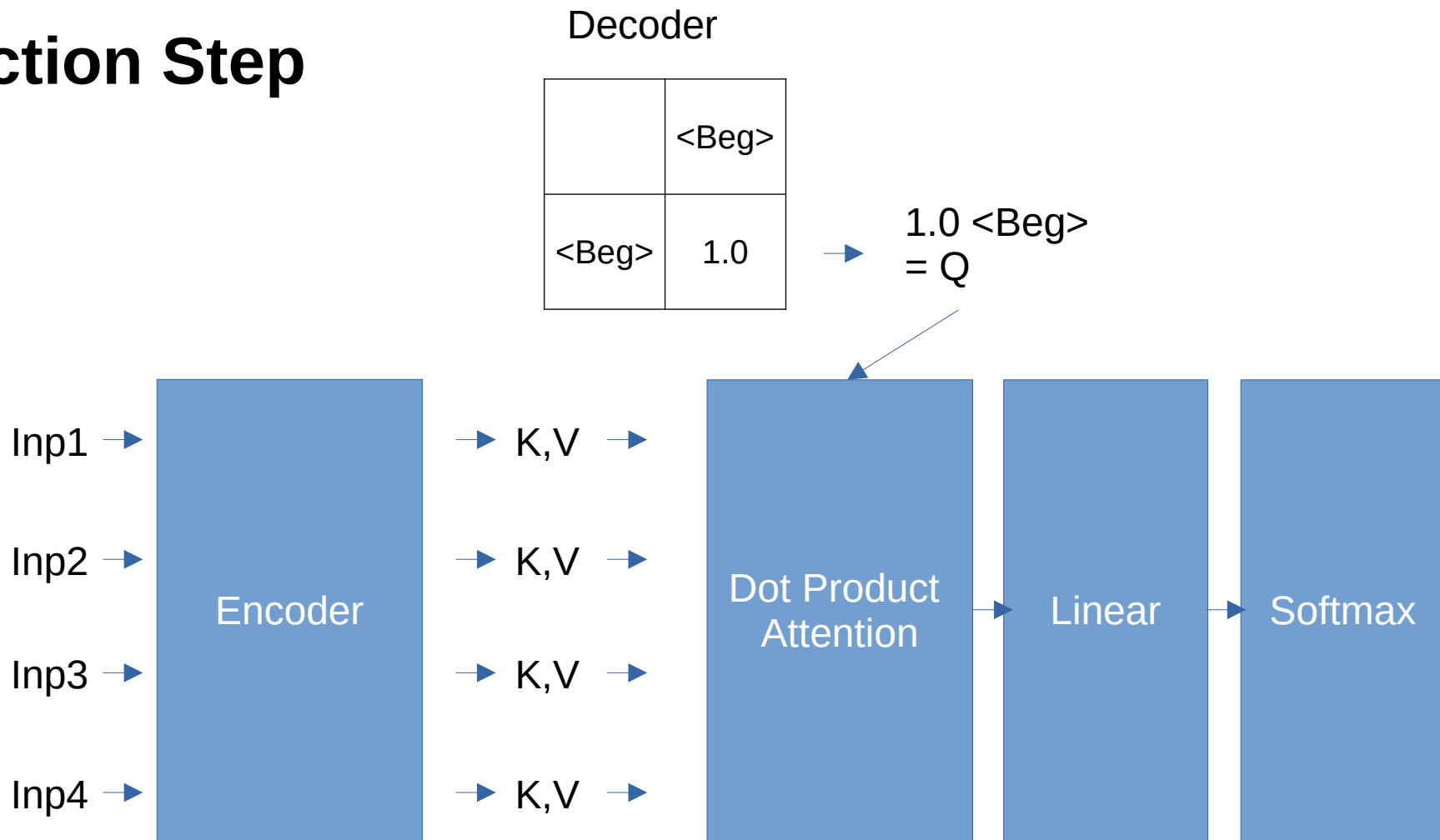
Prediction Step



Decoder

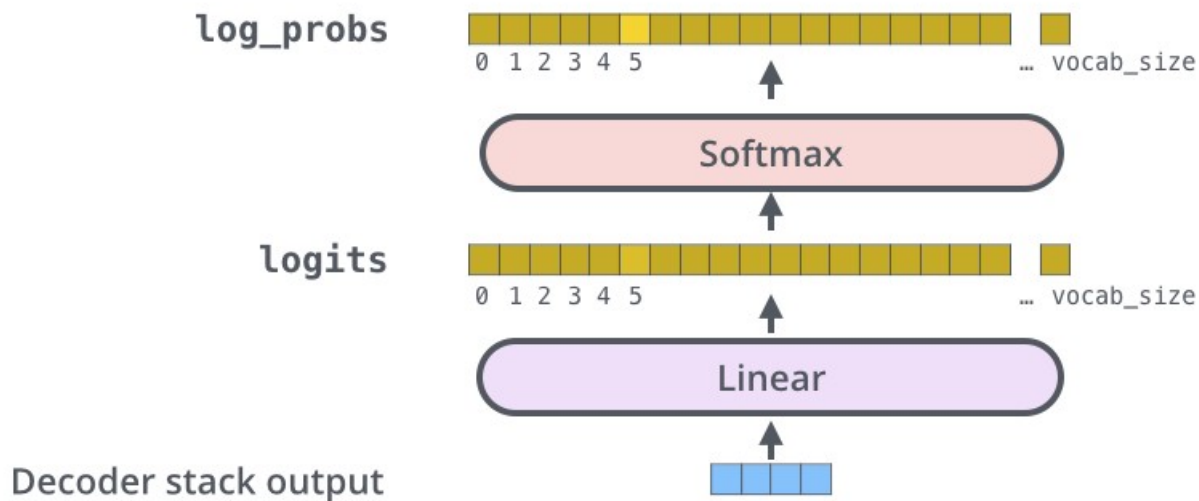


Prediction Step



Which word in our vocabulary
is associated with this index?

Get the index of the cell
with the highest value
(**argmax**)

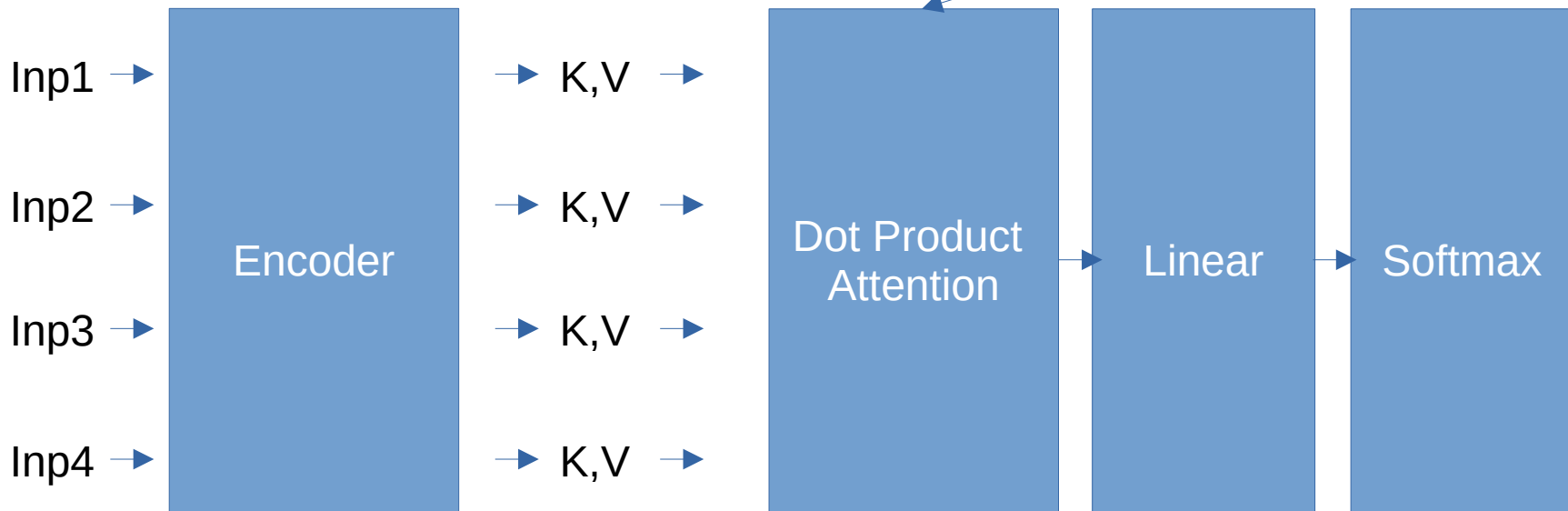


Prediction Step

Decoder

	<Beg>	Out1
Out1	0.3	0.7

$$\begin{aligned} &0.3 \text{ <Beg>} \\ &+ 0.7 \text{ Out1} \\ &= Q \end{aligned}$$

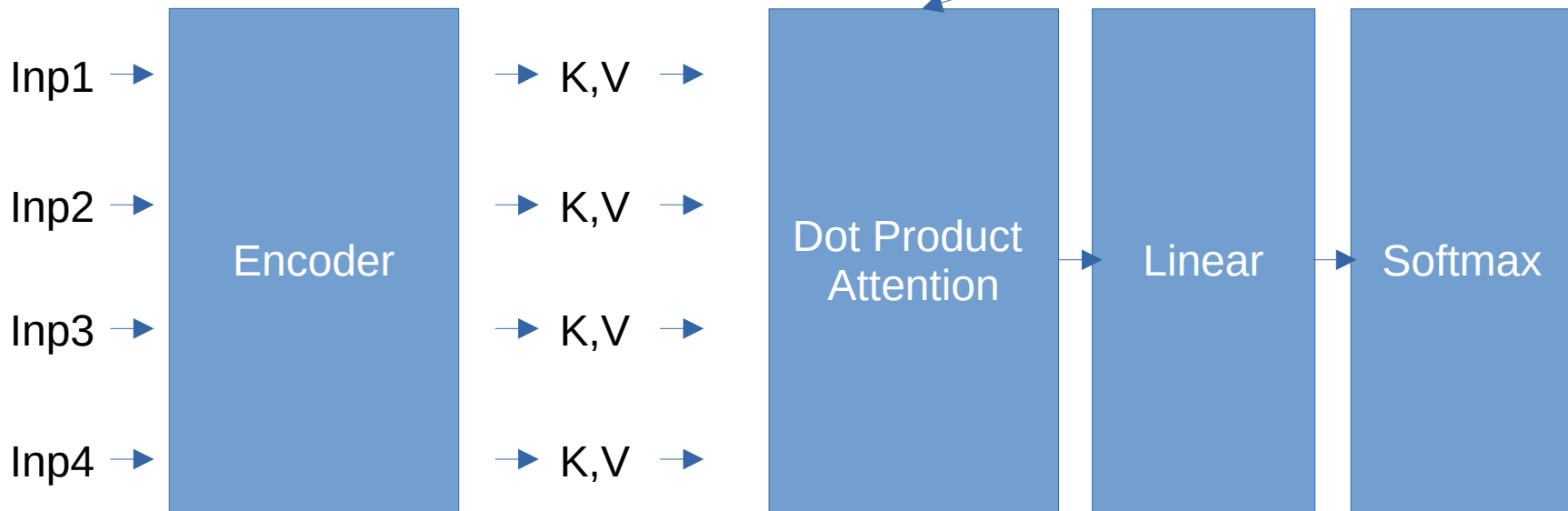


Prediction Step

Decoder

	<Beg>	Out1	Out2
Out2	0.3	0.1	0.6

$$\begin{aligned} &0.3 \text{ <Beg>} \\ &+ 0.1 \text{ Out1} \\ &+ 0.6 \text{ Out2} \\ &= Q \end{aligned}$$



Training Optimization

- What we saw before
 - Iterative prediction
- During training we know the target sentence
 - We can make all predictions in parallel

Masked Self Attention

	Input1	Input2	Input3	Input4	
Input1	→ K,V
Input2	→ K,V
Input3	→ K,V
Input3	→ K,V

	<Beg>	Out1	Out2	Out3	
<Beg>	...				→ Q1
Out1			→ Q2
Out2		→ Q3
Out3	→ Q4

-> parallel prediction of all tokens (training only)