

INF721

2024/2



Deep Learning

L17: Transformers

Logistics

Last Lecture

- ▶ Machine Translation
- ▶ Decoding
 - ▶ Greedy Search
 - ▶ Beam Search
- ▶ Attention in RNNs

Lecture Outline

- ▶ Machine Translation
- ▶ Problems with RNNs
- ▶ Transformers
 - ▶ Self-Attention
 - ▶ Multi-head Attention
 - ▶ Encoder & Decoder
 - ▶ Positional Encoding
 - ▶ Masked Multi-head Attention

Machine Translation

Given a dataset of sentence pairs:

$$(x = \{x^{<1>}, x^{<2>}, \dots, x^{<T_x>}\}, y = \{y^{<1>}, y^{<2>}, \dots, y^{<T_y>}\}),$$

we want to learn a model that maps x into y .

Portuguese

English

Olá, como vai você?

Hello, how are you?

O livro está em cima da mesa.

The book is on the table.

Lucas irá viajar ao Rio em Dezembro.

Lucas is travelling to Rio in December.

Em Dezembro, Lucas irá viajar ao Rio.

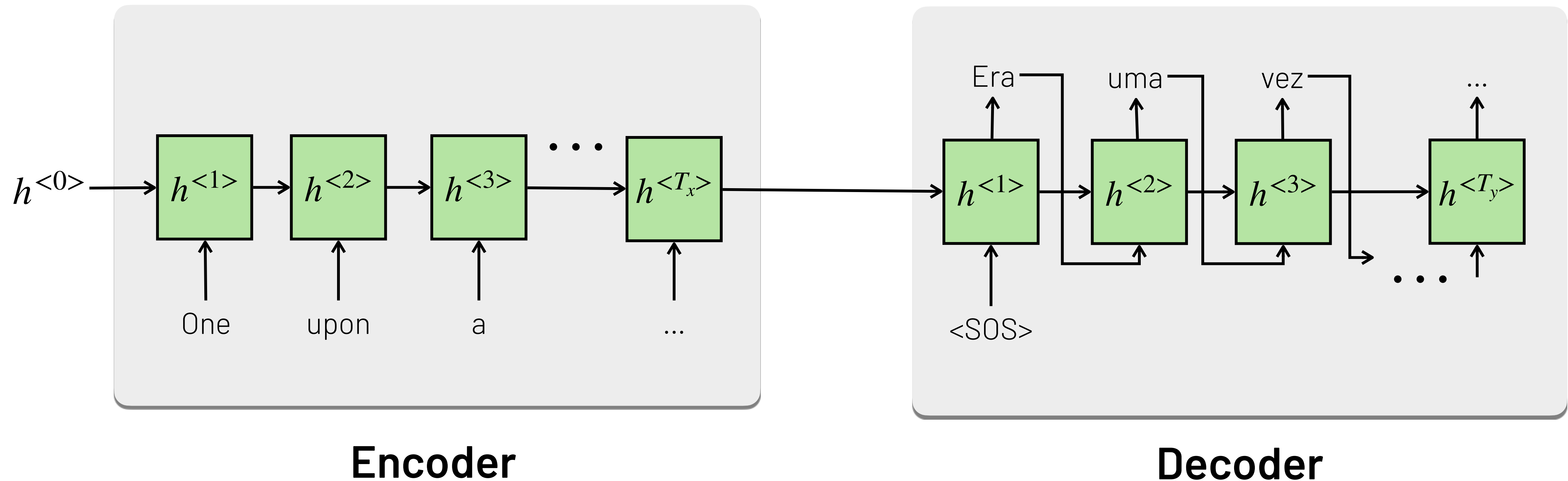
Lucas is travelling to Rio in December.

....

....

Problems with RNNs

- ▶ Struggle to capture long dependencies in sequences
- ▶ Hard to parallelize

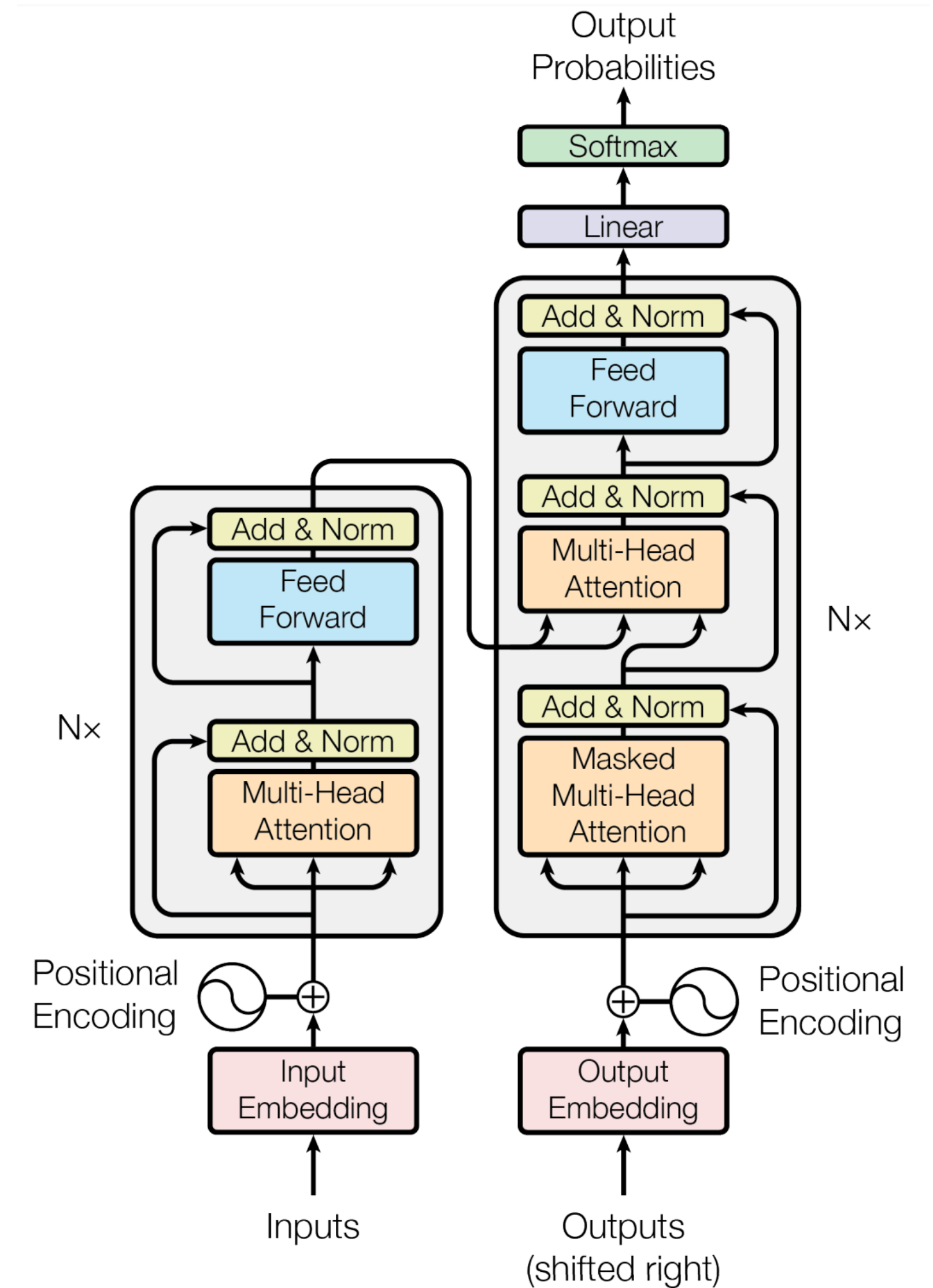


Transformers

Transformers are an encoder-decoder architecture to process sequences using only attention (eliminating recurrence).

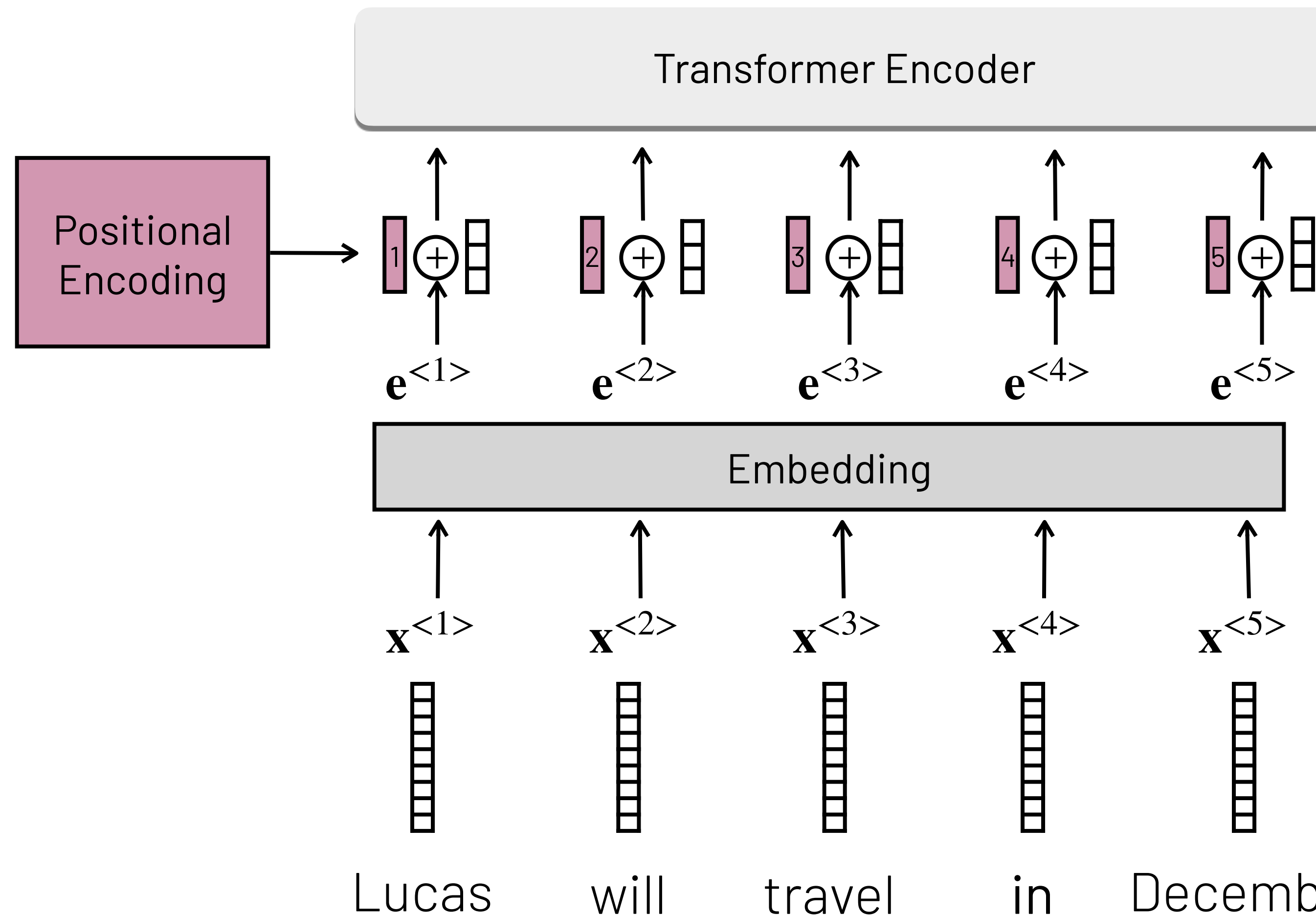
Initially proposed for machine translation, but proved to be very effective in many other problems in:

- ▶ Natural Language Processing
- ▶ Computer Vision
- ▶ Reinforcement Learning
- ▶ ...



Encoder Input

The transformer encoder takes as input a sequence of word embeddings summed with positional encodings. This sequence has the constant size (T_x, d_{model}) throughout the entire model



The contextual embeddings size is typically called d_{model}

$$e_{x+p} = \{e_x^{<1>} + pe^{<1>}, \dots, e_x^{<T_x>} + pe^{<T_x>}\}$$

$(T_x, d_{model}) \rightarrow$ word + positional embedding

$$e_x = \{e_x^{<1>}, \dots, e_x^{<T_x>}\}$$

$(T_x, d_{model}) \rightarrow$ word embedding

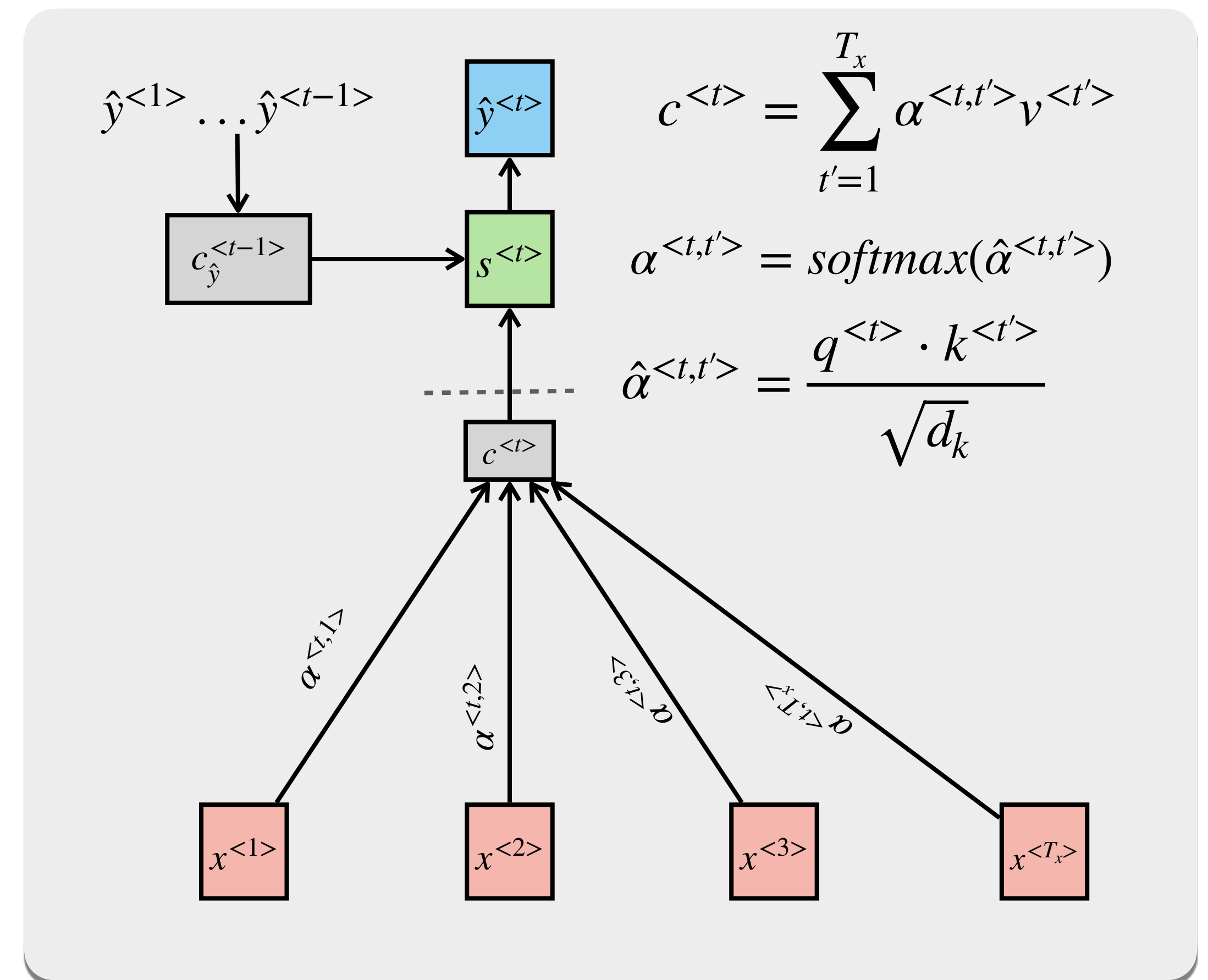
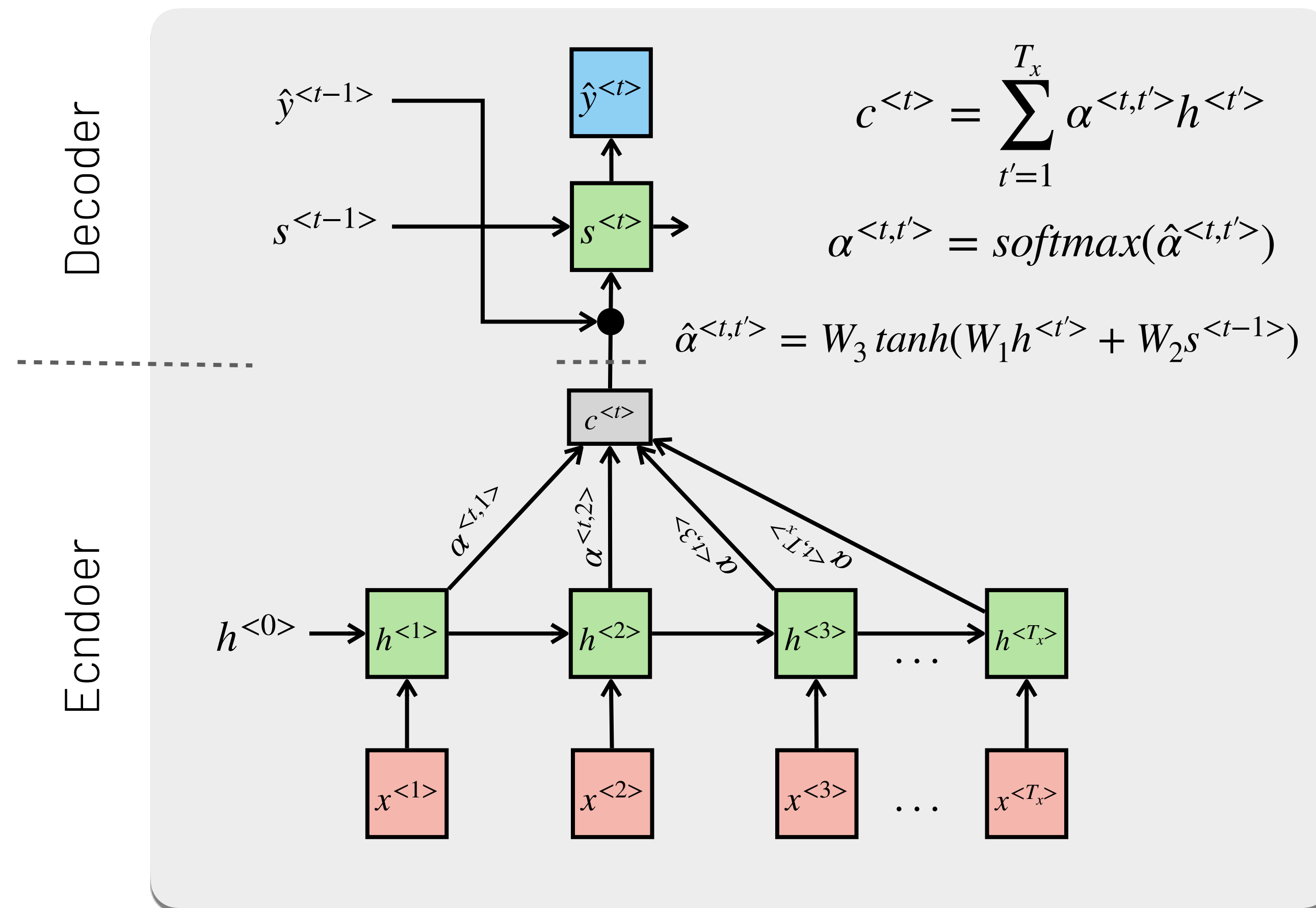
$$x = \{x^{<1>}, \dots, x^{<T_x>}\}$$

$(T_x, |V|) \rightarrow$ 1-hot

Attention in RNNs vs. Transformers

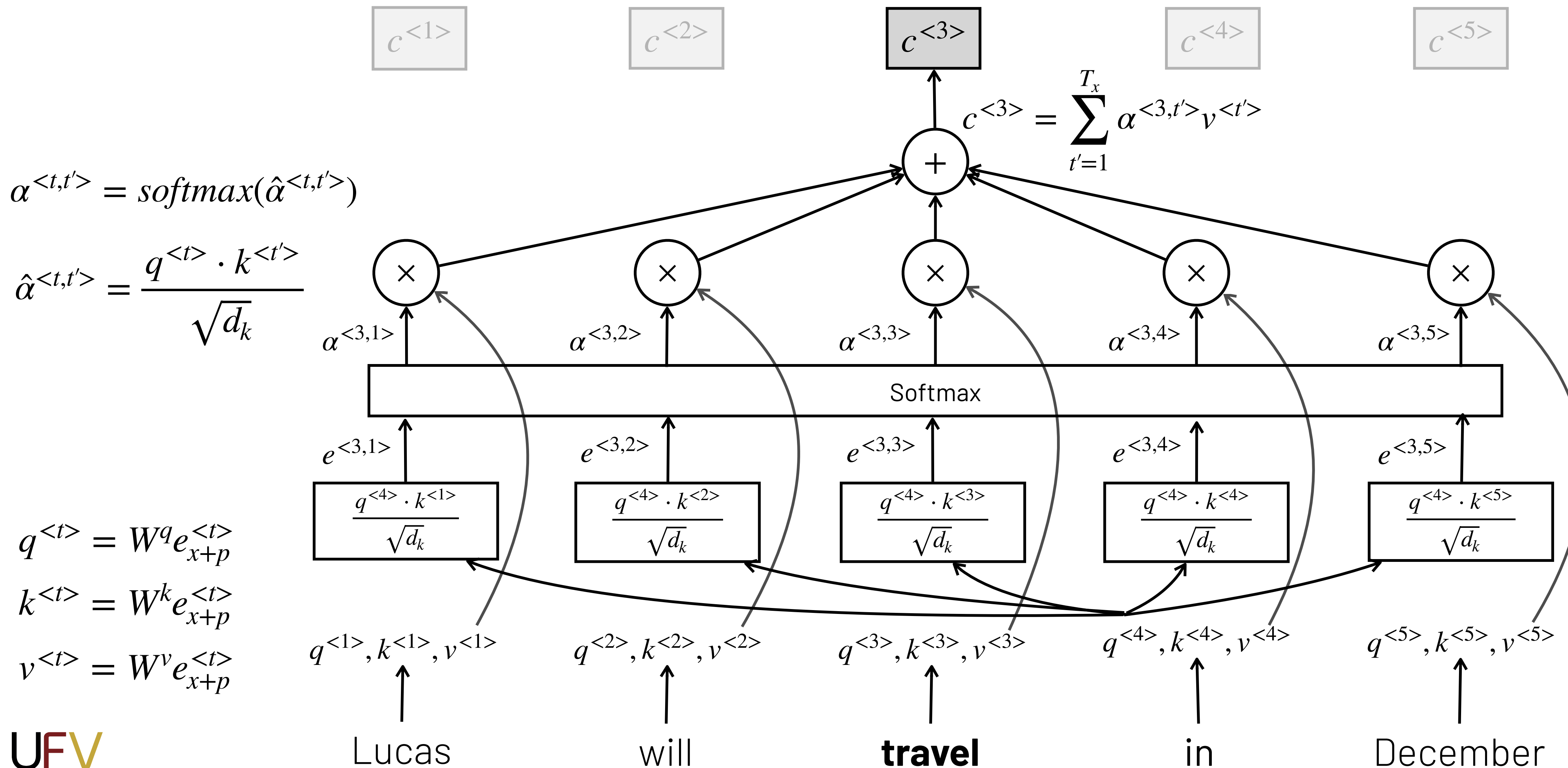
RNNs Bahdanau (Additive) Attention

Transformers Scaled Dot-Product Attention



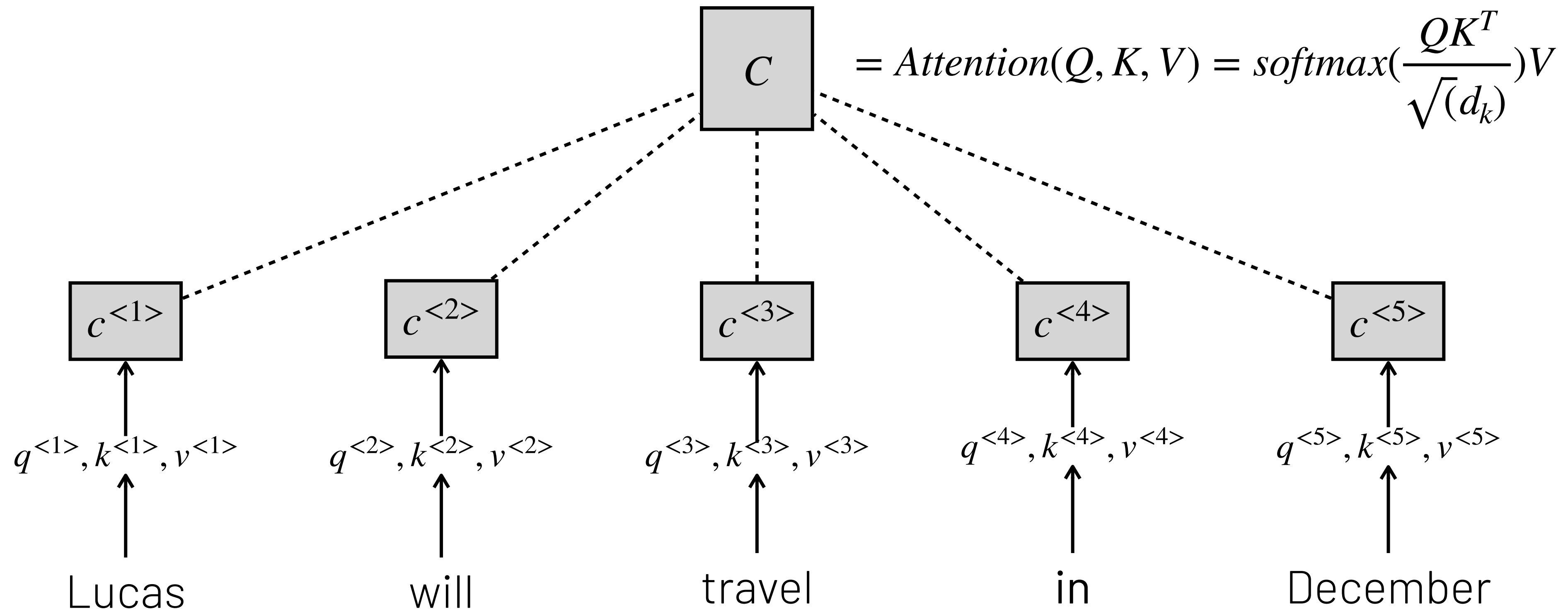
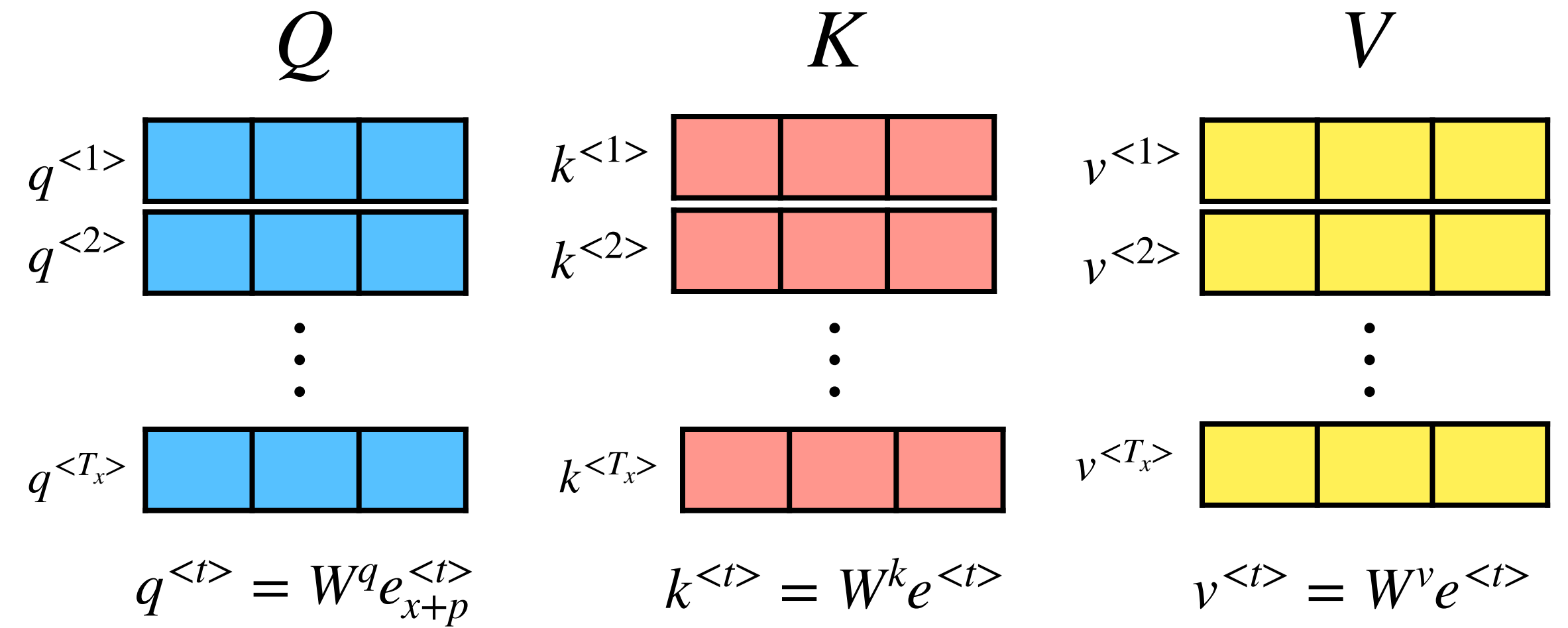
Self-Attention

The key idea behind the Transformer is the **self-attention mechanism**, which learns a context vector $c^{<t>}$ for each input element $x^{<t>}$ based on the input sequence x itself.



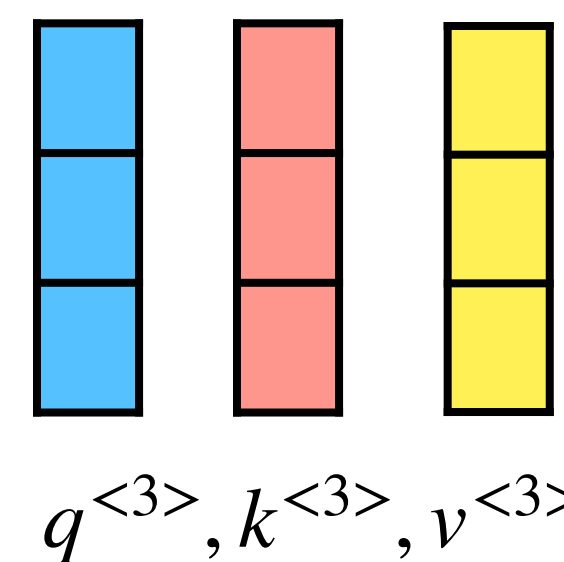
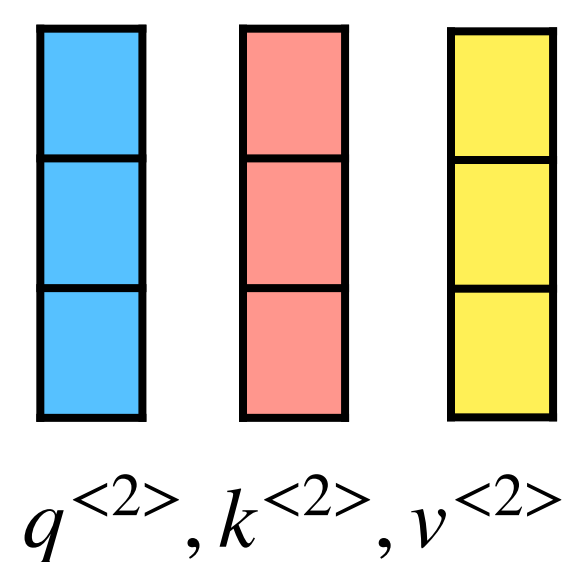
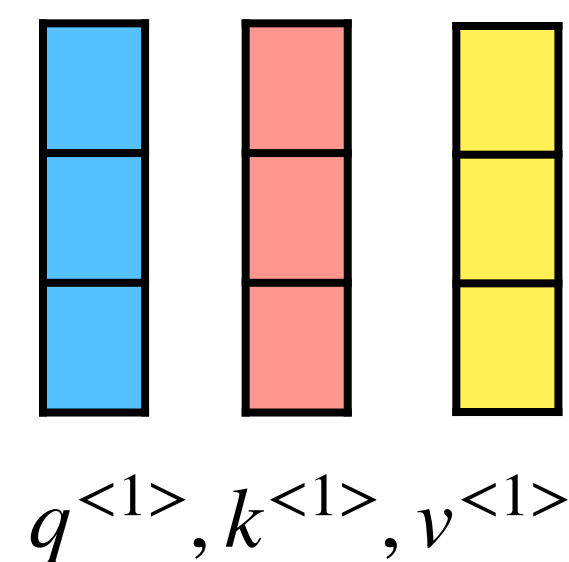
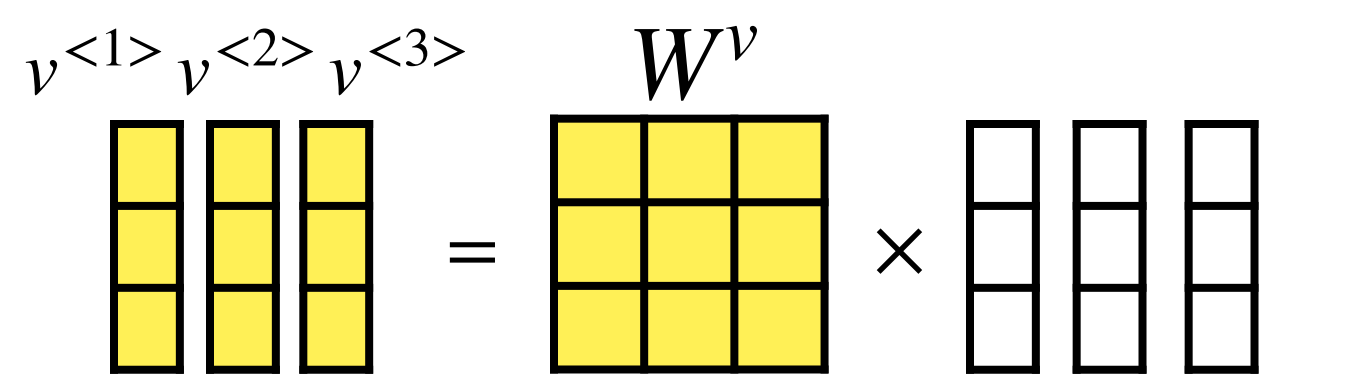
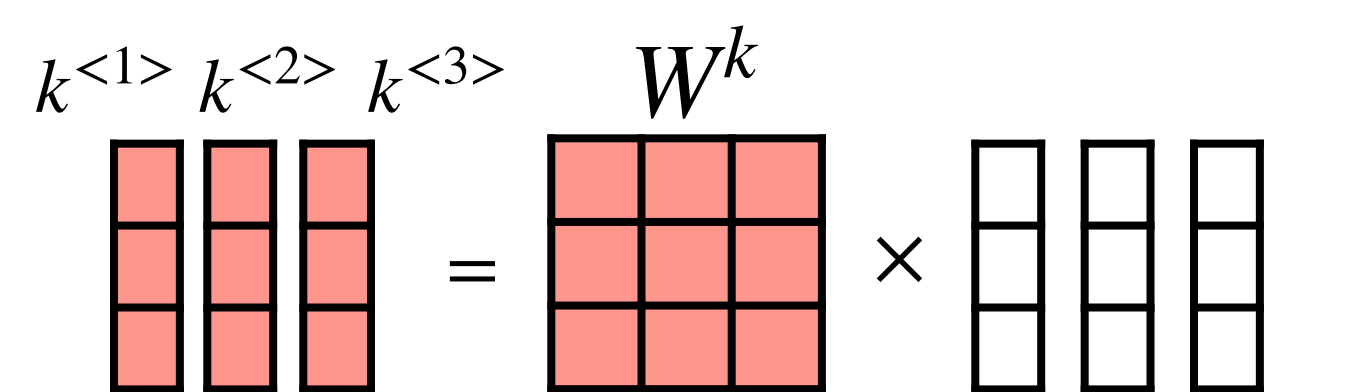
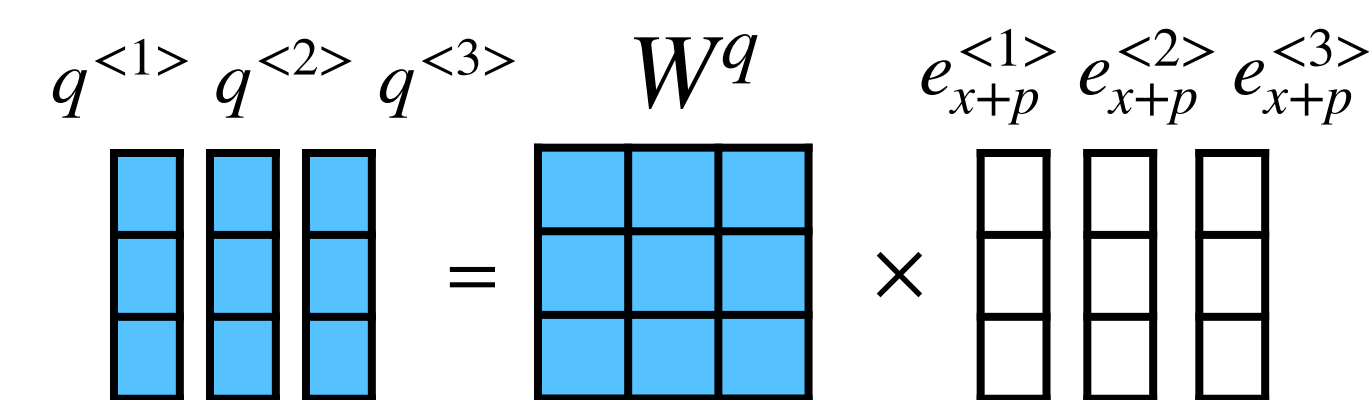
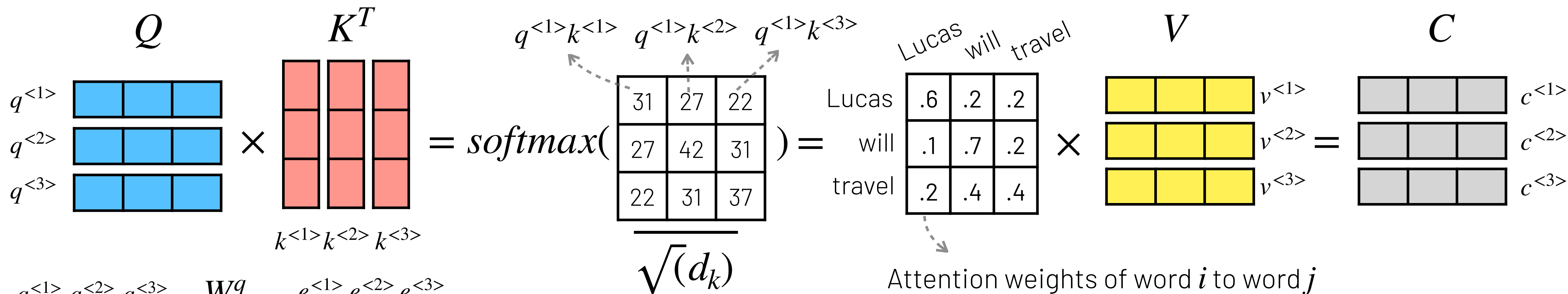
Self-Attention

The contextual representation $C = \{c^{<1>}, \dots, c^{<T_x>}\}$ of the entire input sequence $x = \{x^{<1>}, \dots, x^{<T_x>}\}$ can be computed in a vectorized way combining vectors $q^{<t>}, k^{<t>}, v^{<t>}$ in matrices Q, K e V



Self-Attention

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

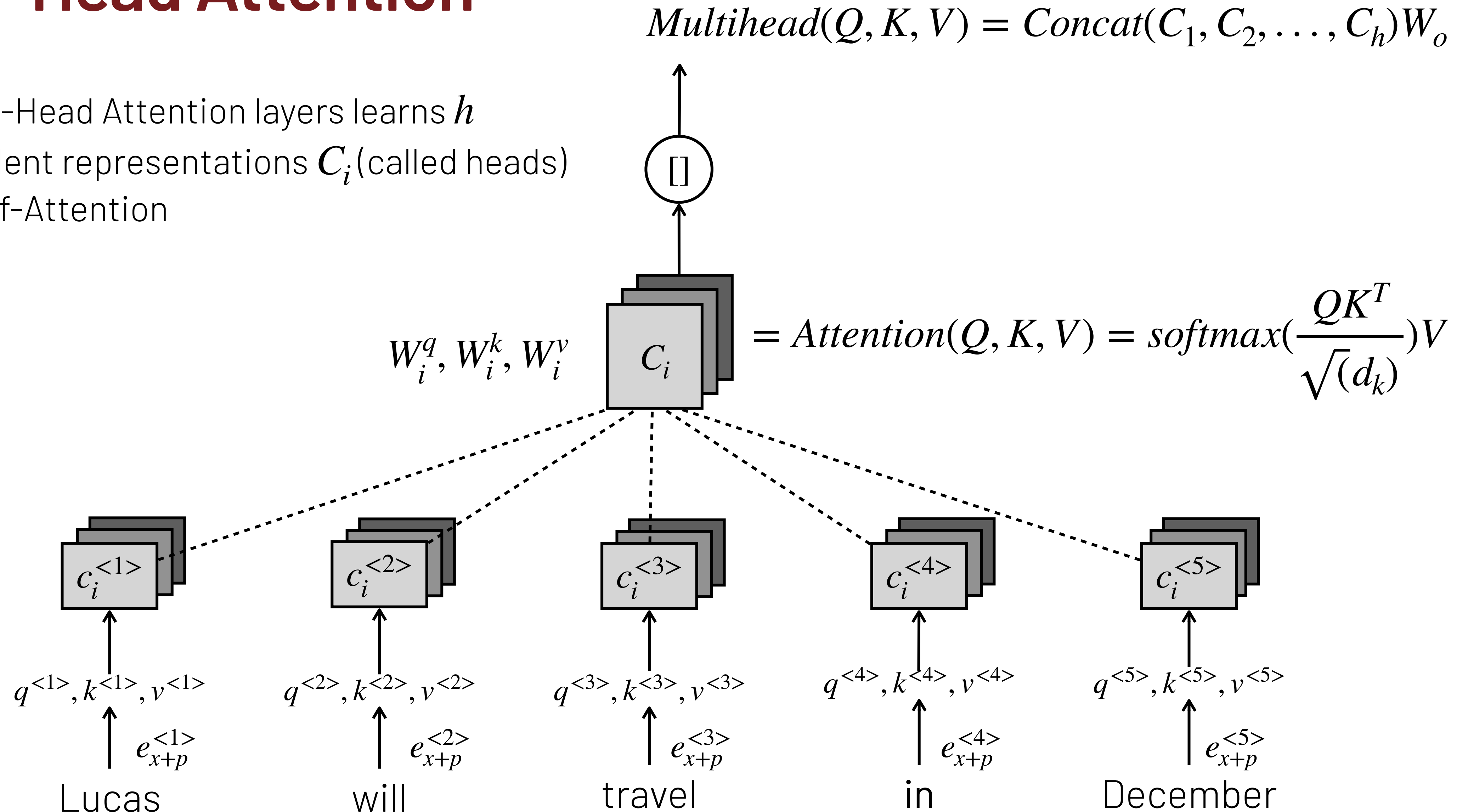


$d_e = d_q = d_k = d_v = 4$
 The sizes of key and query have to be the same. But embeddings and value typically also have same sizes.

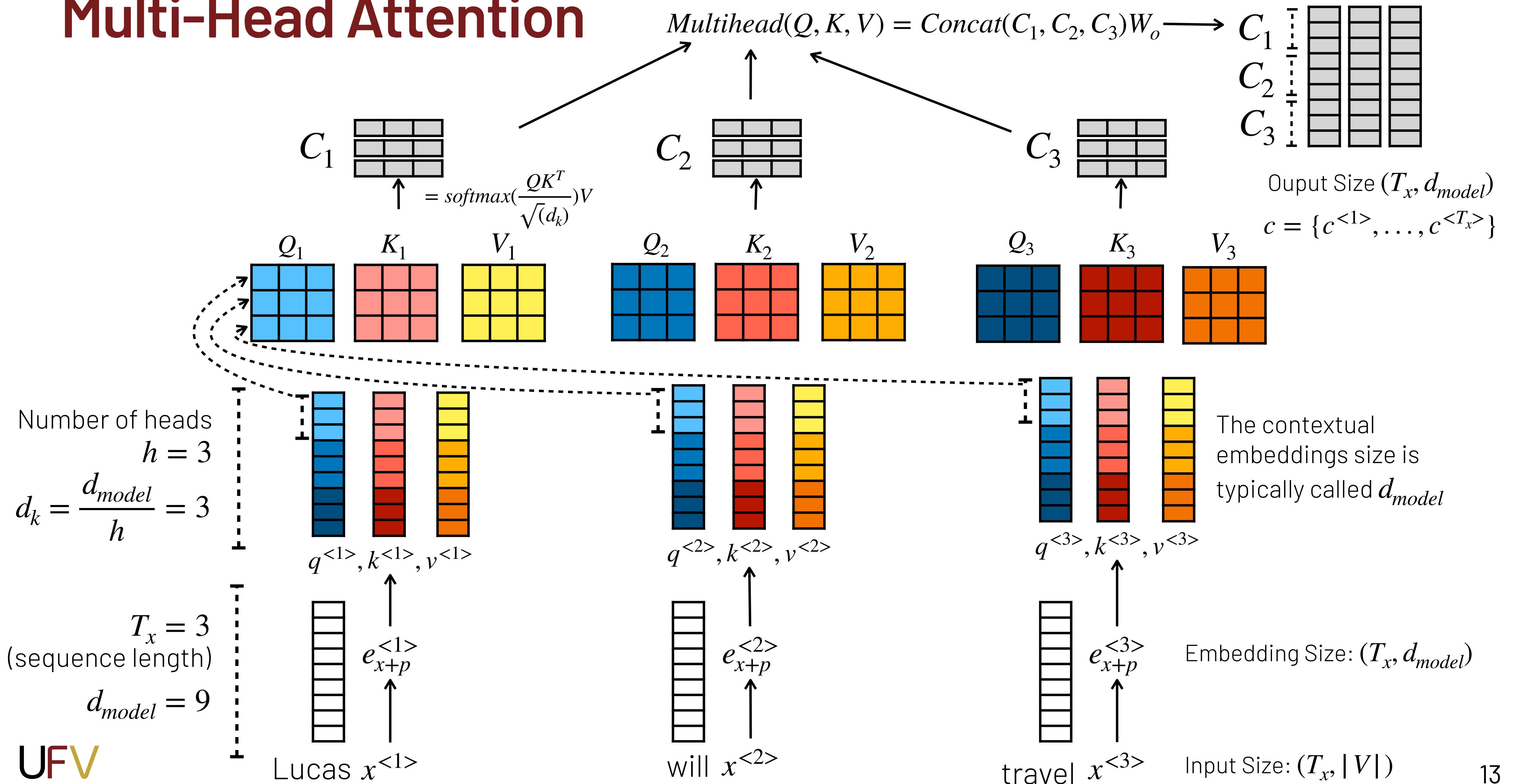
The attention layer receives the embedded vectors $e^{<t'>}$ as input

Multi-Head Attention

The Multi-Head Attention layers learns h independent representations C_i (called heads) using Self-Attention



Multi-Head Attention



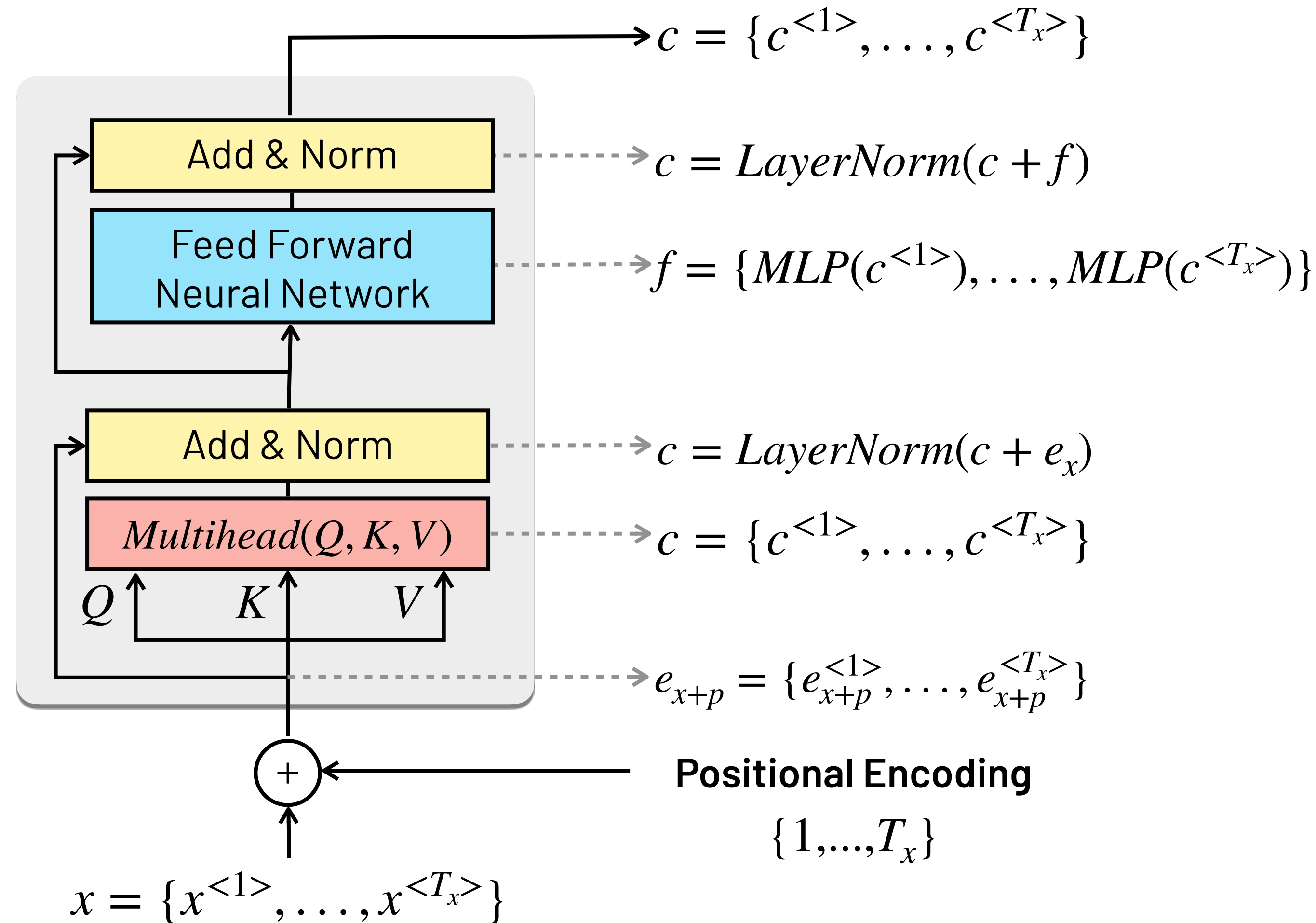
Encoder

Input: a sequence $x = \{x^{<1>}, \dots, x^{<T_x>}\}$

Output: a contextual representation $C = \{c^{<1>}, \dots, c^{<T_x>}\}$ of x

The encoder applies a **Multihead Layer** followed by a **Feed Forward Neural Network** (MLP).

Both are normalized with Layer Norm (**Norm**) and connected with a residual connection (**Add**)

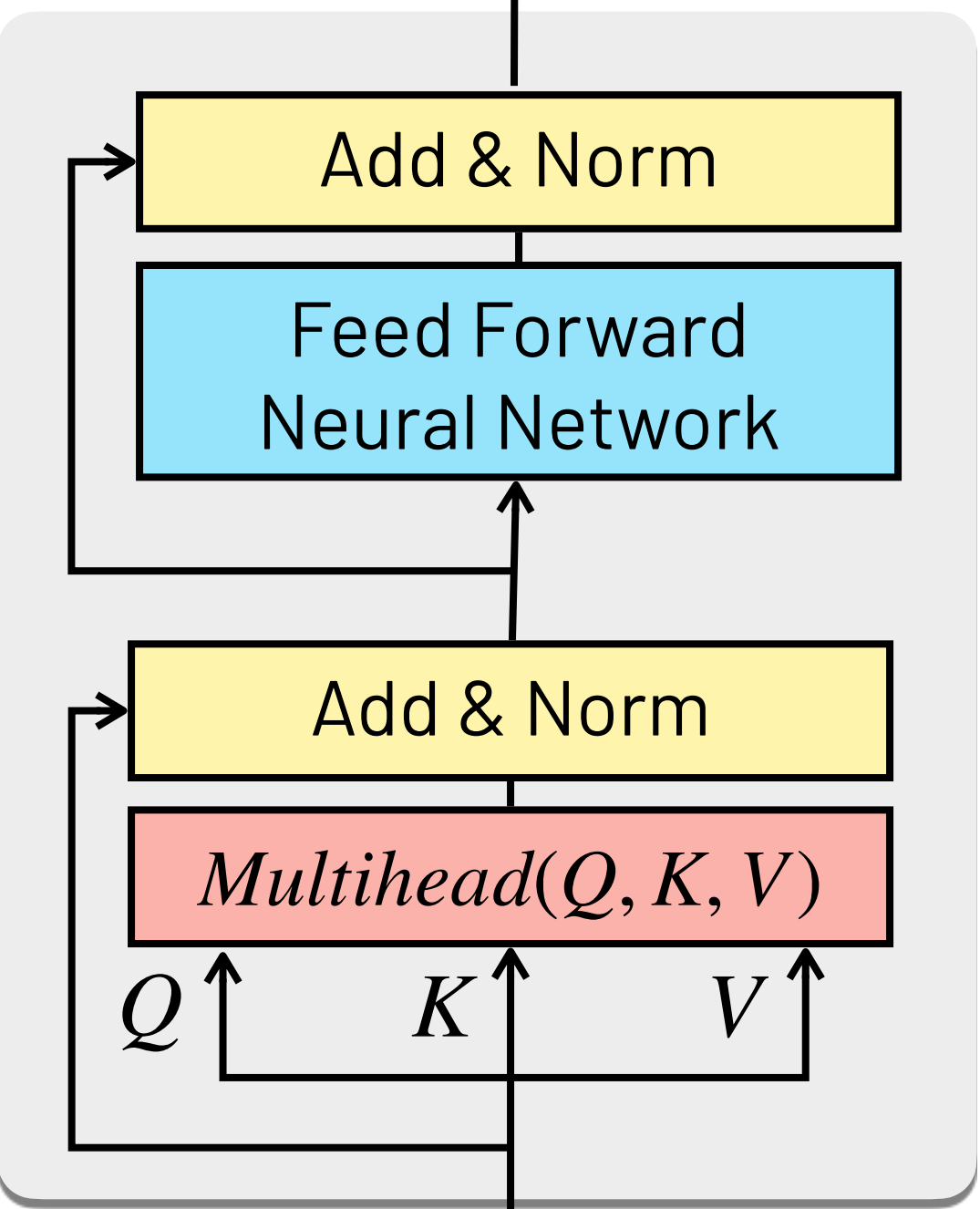


Decoder

Input: Input context C and previous $\{\hat{y}^{<1>}, \dots, \hat{y}^{<t-1>}$ output tokens

Output: The next token $\hat{y}^{<t>}$

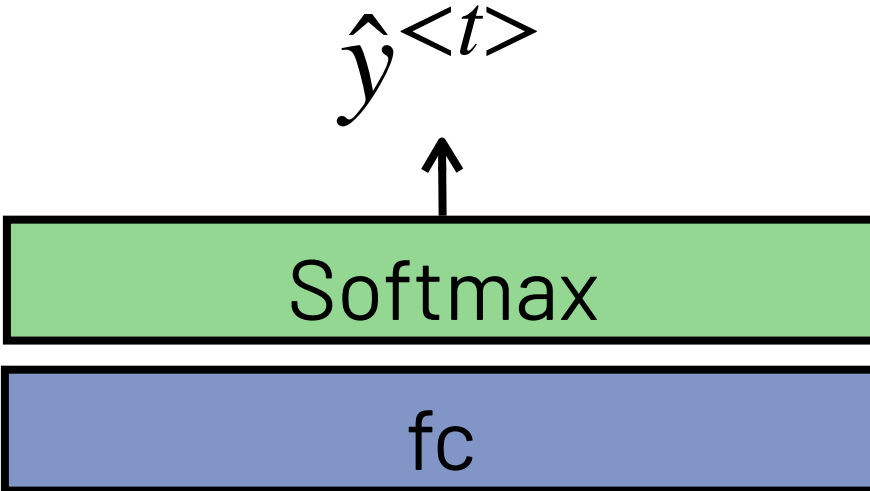
Nx Encoder $c = \{c^{<1>}, \dots, c^{<T_x>}$



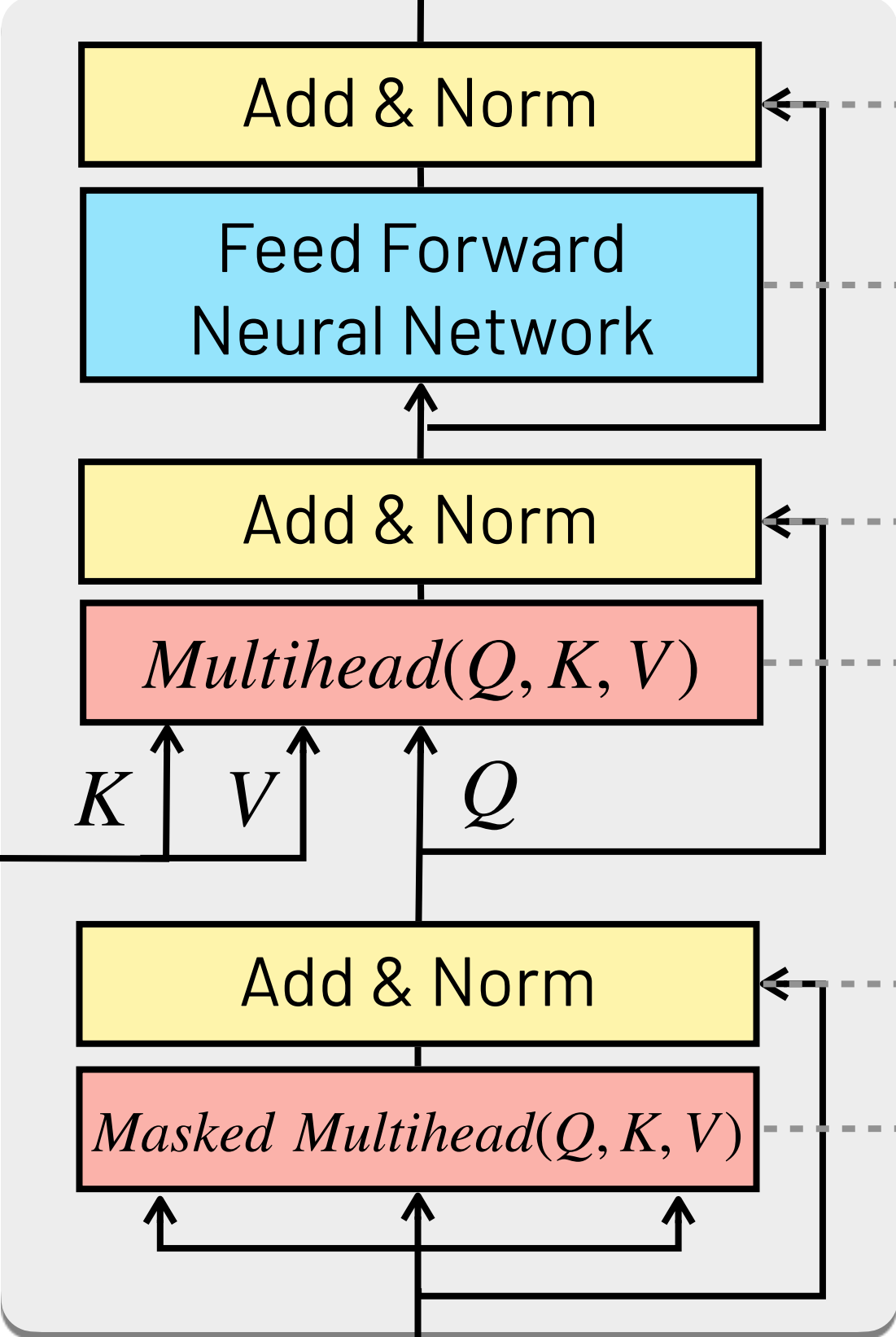
+

$x = \{x^{<1>}, \dots, x^{<T_x>}$

Positional Encoding
 $\{1, \dots, T_x\}$



Decoder Nx



$$K = W^k c$$

$$V = W^v c$$

$K \quad V \quad Q$

$$a_2 = \{a_2^{<1>}, \dots, a_2^{<t-1>}\}$$

$$a_2 = \text{LayerNorm}(a_2 + f)$$

$$f = \{MLP(a_2^{<1>}), \dots, MLP(a_2^{<t-1>})\}$$

$$a_2 = \text{LayerNorm}(a_2 + a_1)$$

$$a_2 = \{a_2^{<1>}, \dots, a_2^{<t-1>}\}$$

$$Q = W^q a_1$$

$$a_1 = \text{LayerNorm}(a_1 + e_{\hat{y}})$$

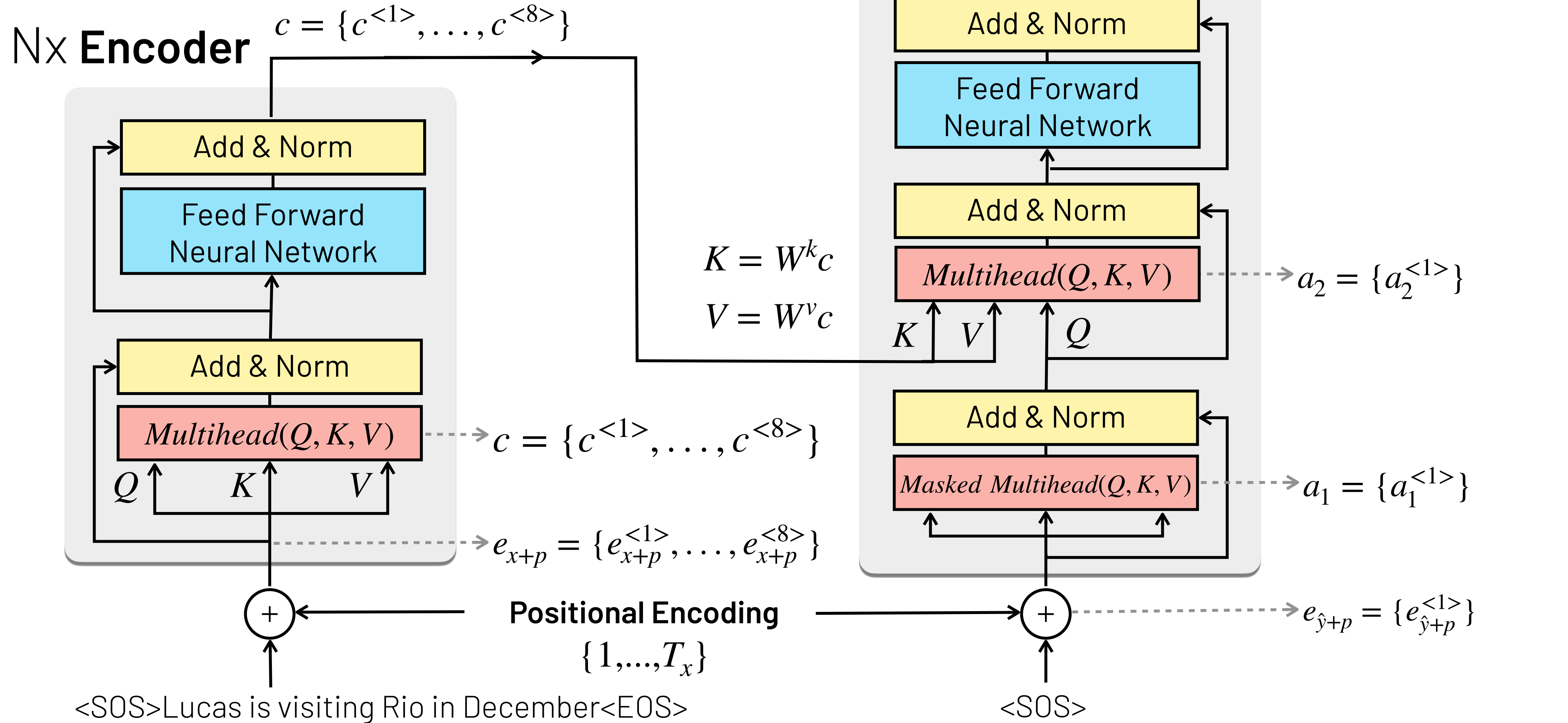
$$a_1 = \{a_1^{<1>}, \dots, a_1^{<t-1>}\}$$

$$e_{\hat{y}} = \{e_{\hat{y}+p}^{<1>}, \dots, e_{\hat{y}+p}^{<t-1>}\}$$

$\hat{y} = \{\hat{y}^{<1>}, \dots, \hat{y}^{<t-1>}\}$

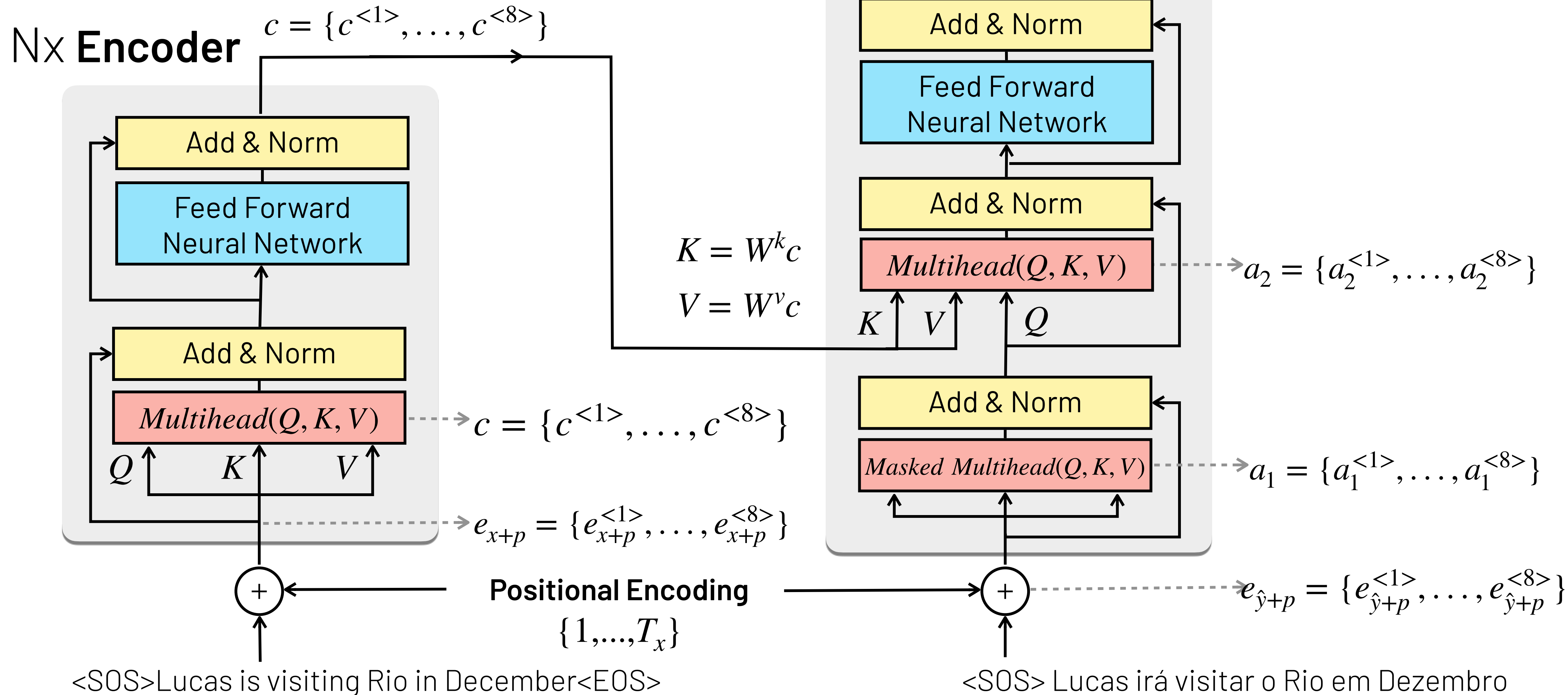
Inference Example

At inference time, the model is autoregressive, i.e., it generates one word at a time.



Inference Example

At inference time, the model is autoregressive, i.e., it generates one word at a time.

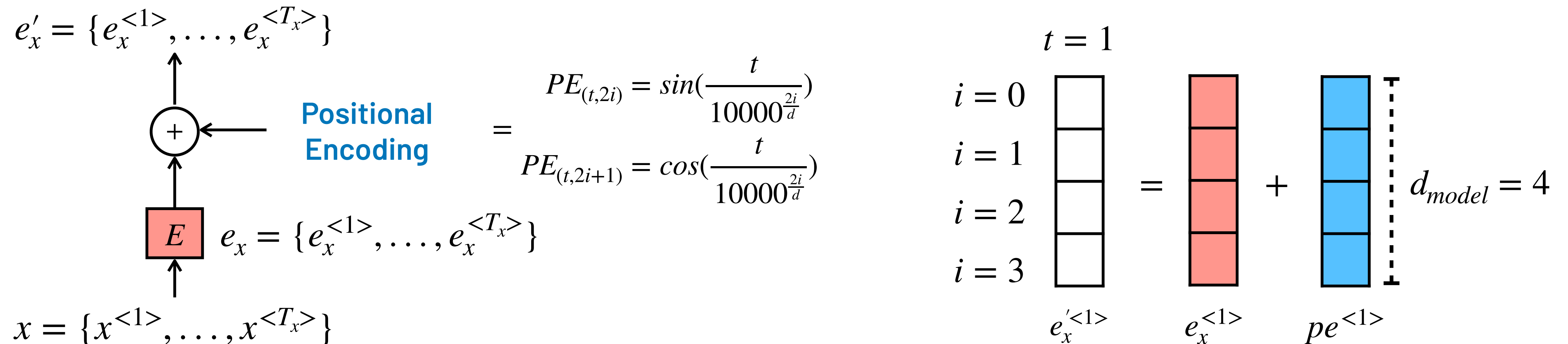


Positional Encoding

The self-attention mechanism does not consider the position of the words.

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

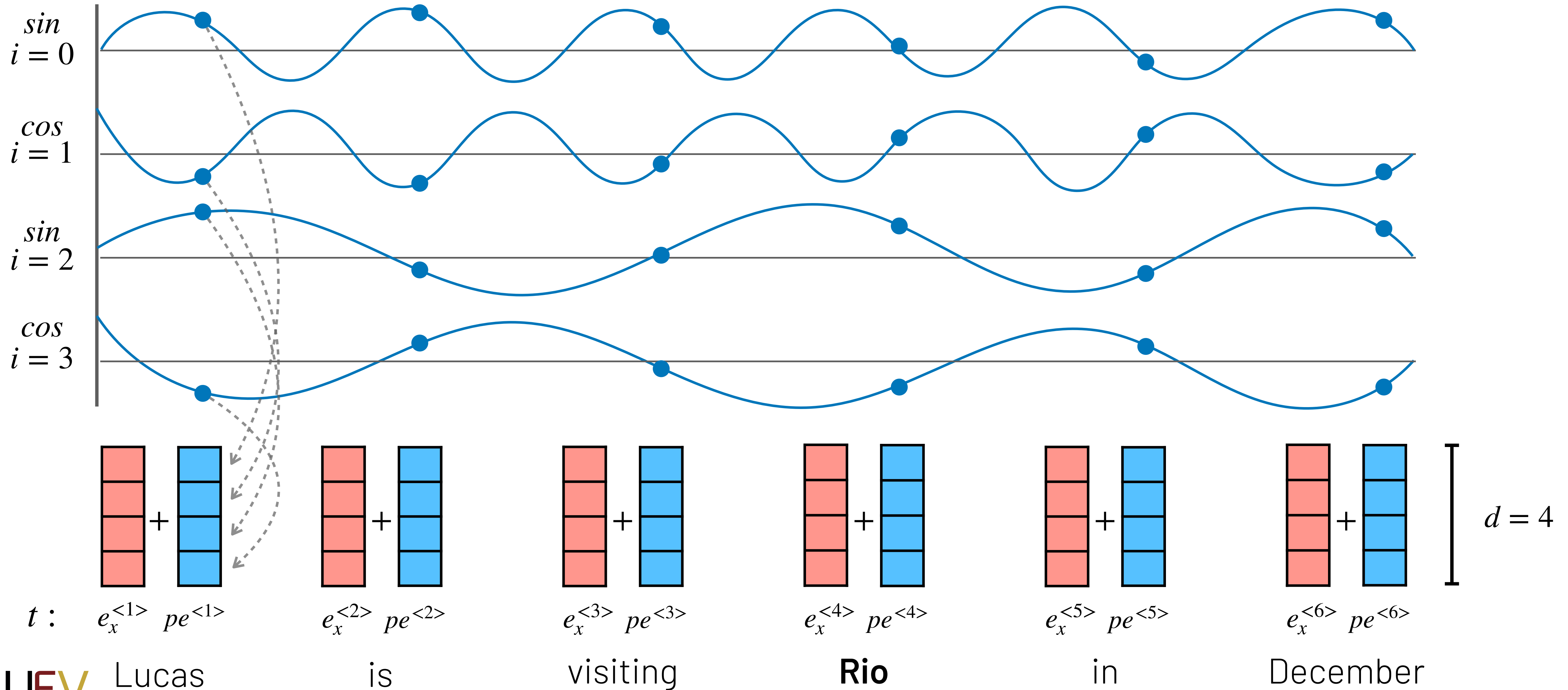
To add this information to the learned contextual representation C , both encoder and decoder add a positional information to each element $x^{<t>}$ of the input $x = \{x^{<1>}, \dots, x^{<T_x>}\}$



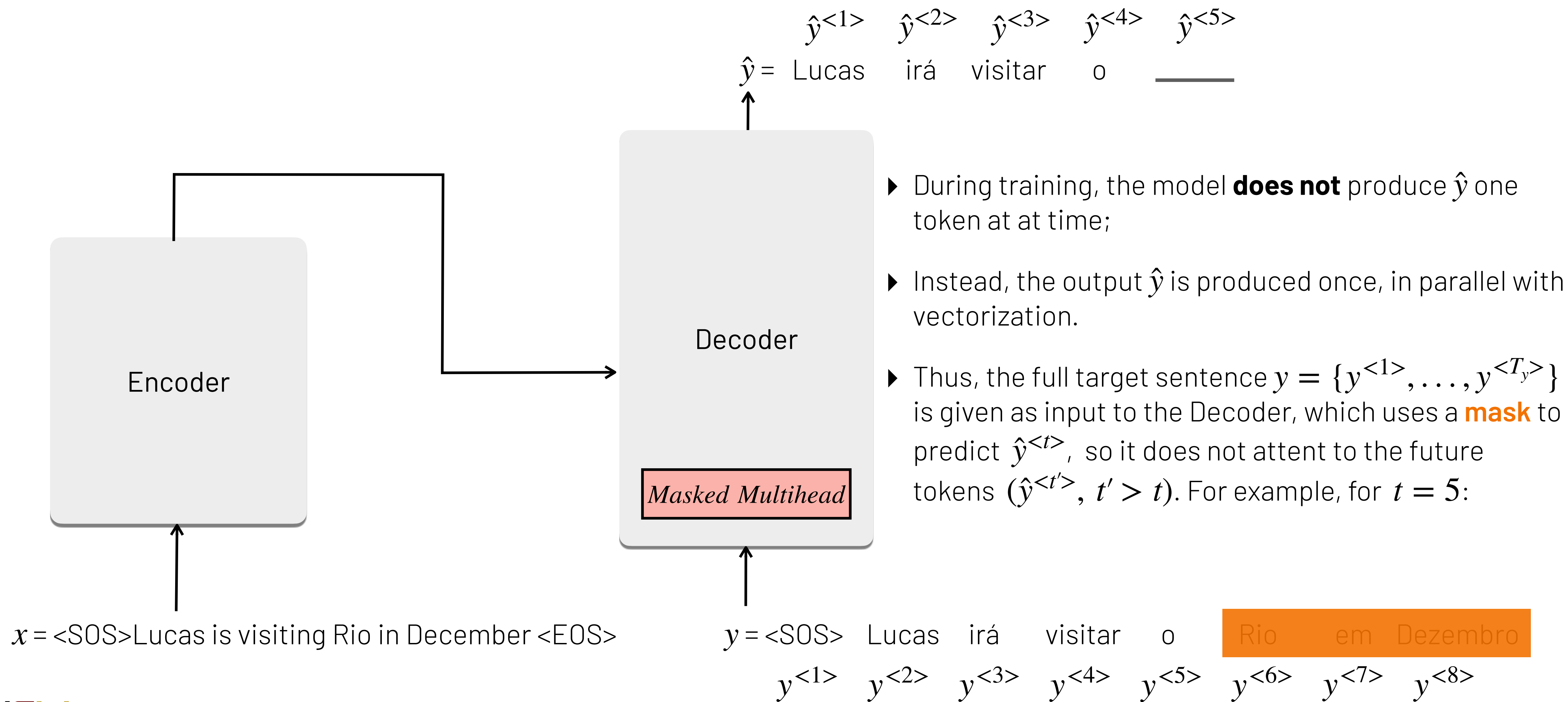
Positional Encoding

$$PE_{(t,2i)} = \sin\left(\frac{t}{10000^{\frac{2i}{d}}}\right) \quad \text{if } i \text{ is even}$$

$$PE_{(t,2i+1)} = \cos\left(\frac{t}{10000^{\frac{2i}{d}}}\right) \quad \text{if } i \text{ is odd}$$

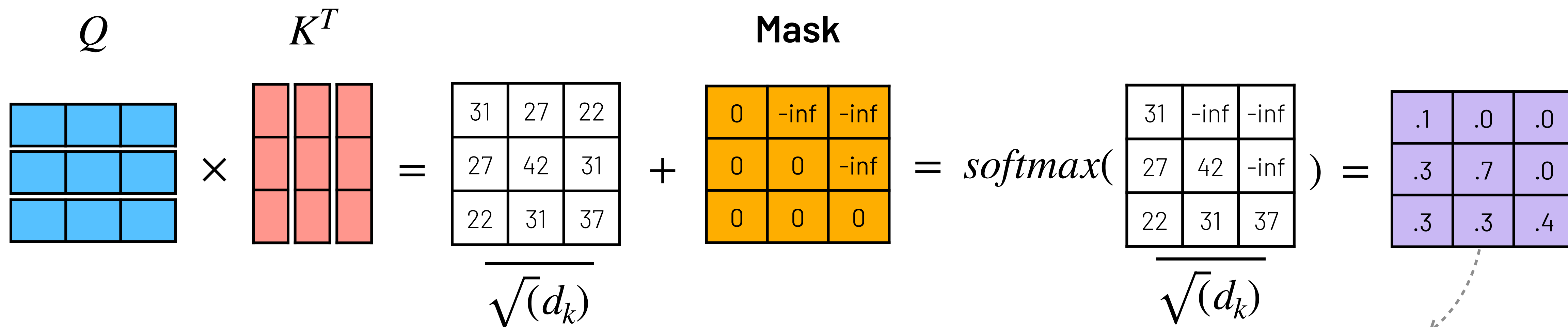


Training with Masked Multi-head Attention



- ▶ During training, the model **does not** produce \hat{y} one token at a time;
- ▶ Instead, the output \hat{y} is produced once, in parallel with vectorization.
- ▶ Thus, the full target sentence $y = \{y^{<1>}, \dots, y^{<T_y>}\}$ is given as input to the Decoder, which uses a **mask** to predict $\hat{y}^{<t>}$, so it does not attend to the future tokens ($\hat{y}^{<t'>}, t' > t$). For example, for $t = 5$:

Training with Masked Multi-head Attention

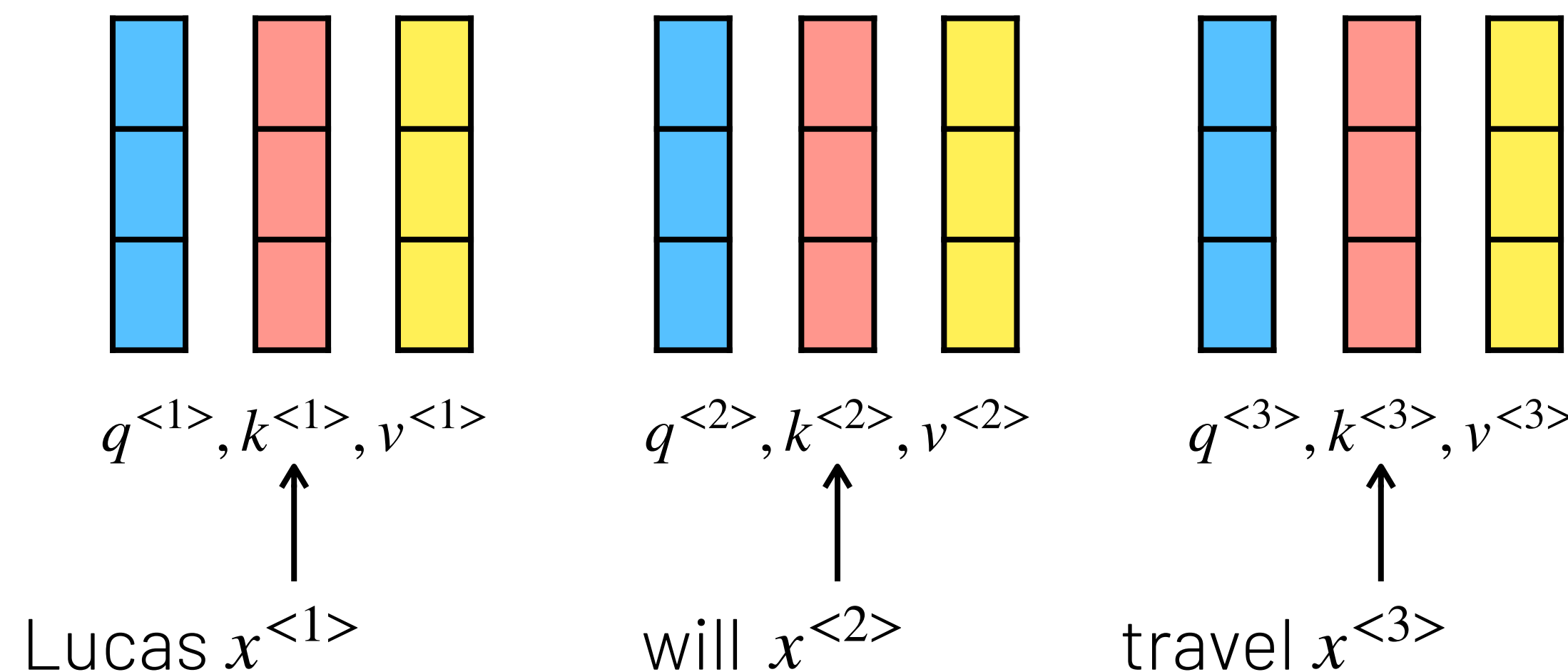


Attention weights (normalized by row) of work i to word j after mask application

$$q^{<t'>} = W^q e^{<t'>}$$

$$k^{<t'>} = W^k e^{<t'>}$$

$$v^{<t'>} = W^v e^{<t'>}$$



Next Lecture

L17: Transformers (Part II)

Case studies of transformers: BERT and GPT