# INF721

2024/2

# Deep Learning

## L17: Transformers

# Logistics

**Last Lecture**

▸ Machine Translation

▸ Decoding

    ▸ Greedy Search

    ▸ Beam Search

▸ Attention in RNNs

UFV

# Lecture Outline

▶ Machine Translation

▶ Problems with RNNs

▶ Transformers

  ▶ Self-Attention

  ▶ Multi-head Attention

  ▶ Encoder & Decoder

  ▶ Positional Encoding

  ▶ Masked Multi-head Attention

UF V
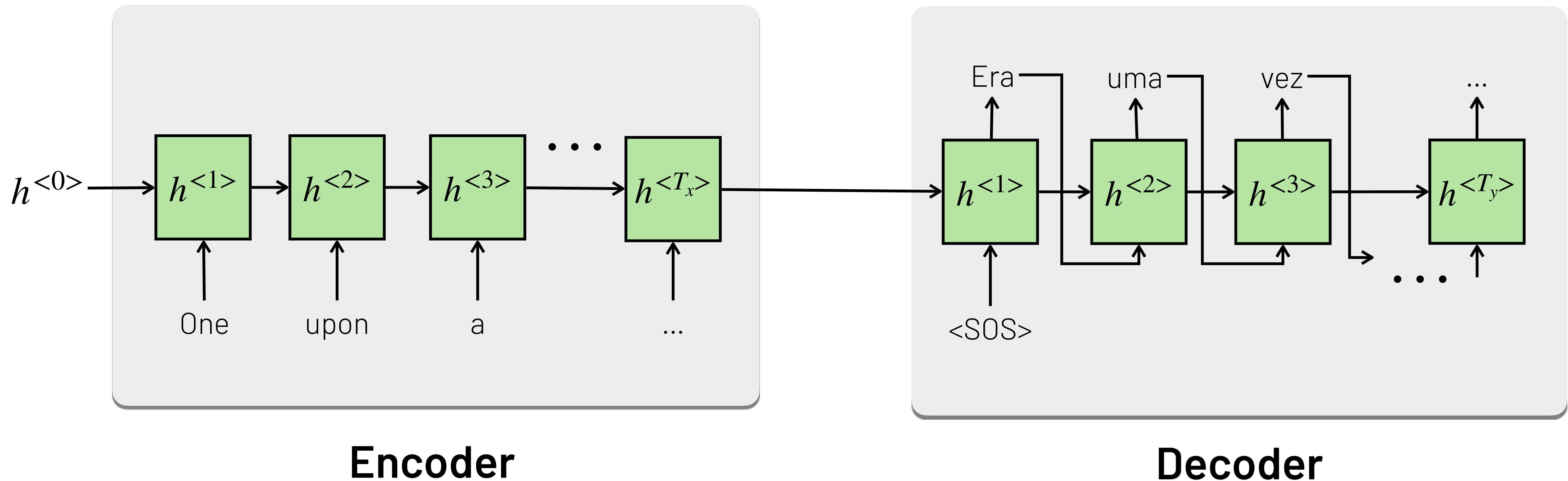
# Machine Translation

Given a dataset of sentence pairs:

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

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. |
| .... | .... |

UFV

# Problems with RNNs

▸ Struggle to capture long dependencies in sequences

▸ Hard to parallelize
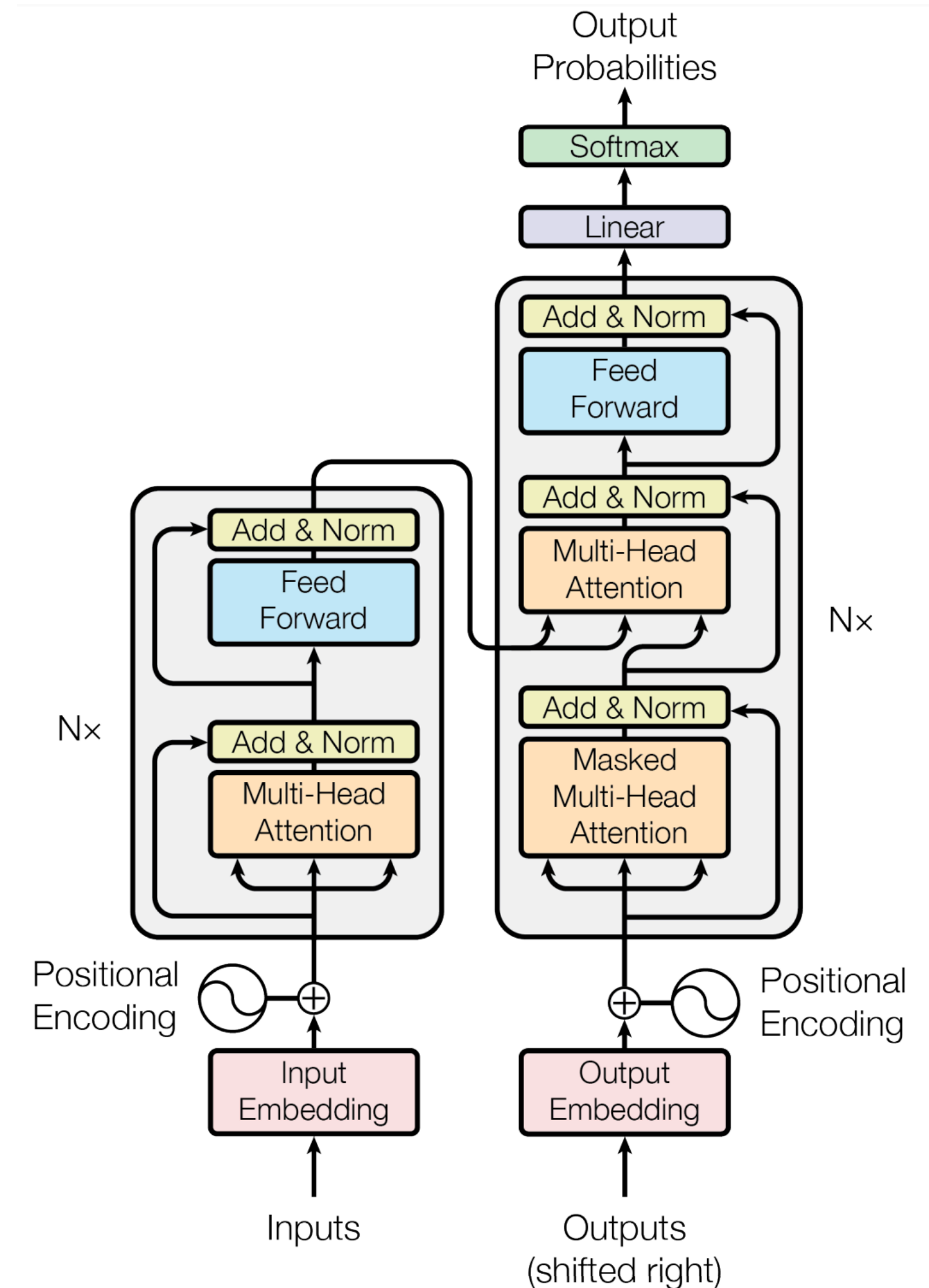


**Encoder**

**Decoder**

# Transformers

**Transformers** are an encoder-decoder architecture to process sequences using only attention (eliminating recurrence).
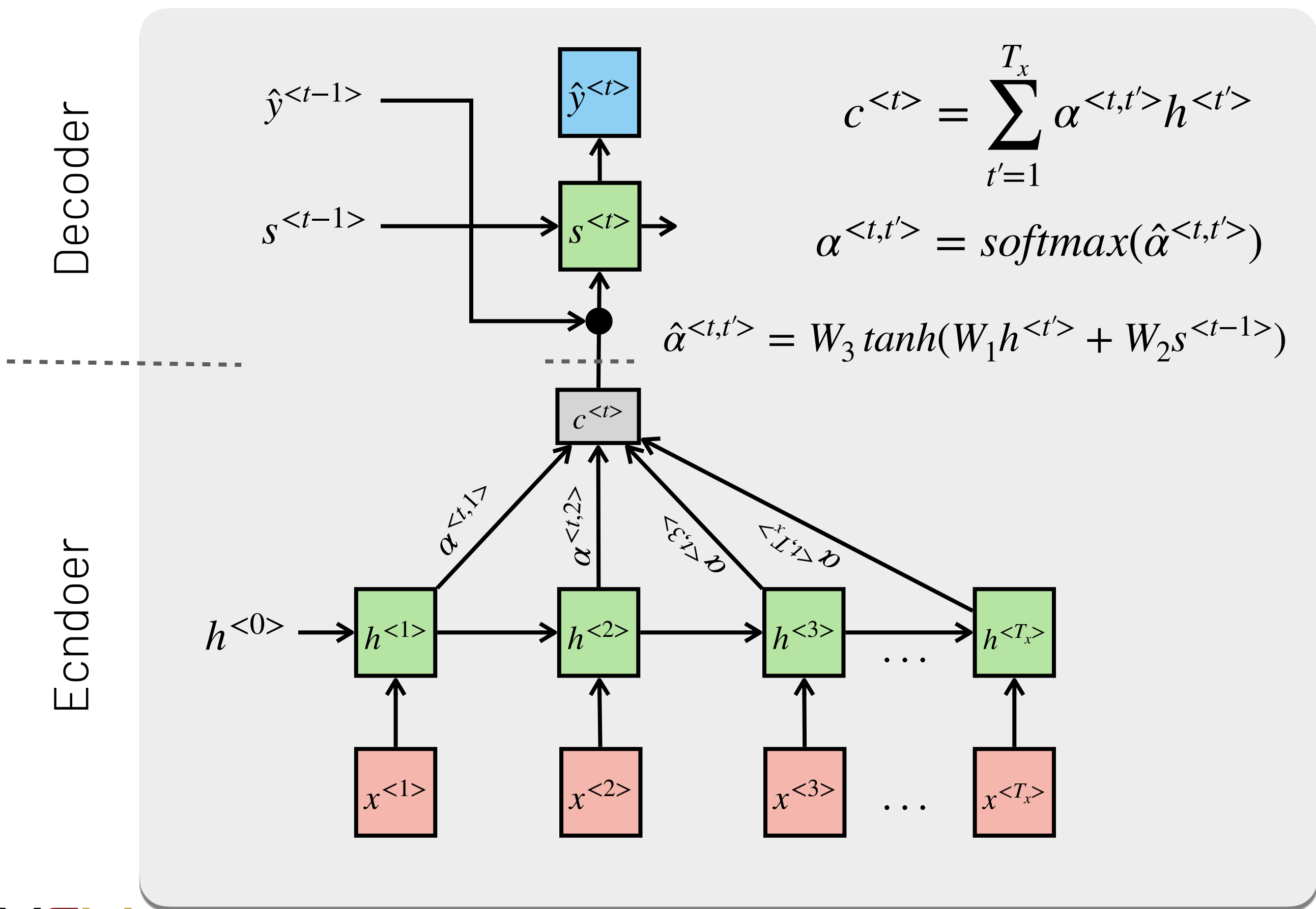
Initially proposed for machine translation, but proved to be very effective in many other problems in:

▸ Natural Language Processing

▸ Computer Vision

▸ Reinforcement Learning

▸ ...



Vaswani et al. 2017, Attention Is All You Need

# Attention in RNNs vs. Transformers

**RNNs**  Badahnau (Additive) Attention

**Transformers**  Scaled Dot-Product Attention



$$c^{<t>} = \sum_{t'=1}^{T_x} \alpha^{<t,t'>} h^{<t'>}$$

$$\alpha^{<t,t'>} = softmax(\hat{\alpha}^{<t,t'>})$$

$$\hat{\alpha}^{<t,t'>} = W_3\, tanh(W_1 h^{<t'>} + W_2 s^{<t-1>})$$

$$c^{<t>} = \sum_{t'=1}^{T_x} \alpha^{<t,t'>} v^{<t'>}$$

$$\alpha^{<t,t'>} = softmax(\hat{\alpha}^{<t,t'>})$$

$$\hat{\alpha}^{<t,t'>} = \frac{q^{<t>} \cdot k^{<t'>}}{\sqrt{d_k}}$$

UFV

# 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 $d_{model}$

$$X_{epe} = X_e + PE$$

$$X_{epe} \in \mathbb{R}^{T_x \times d_{model}} \rightarrow \text{word + pos. embed.}$$

$$X_e = EX$$

$$X_e \in \mathbb{R}^{T_x \times d_{model}} \rightarrow \text{word embedding}$$

$$E \in \mathbb{R}^{|V| \times d_{model}} \rightarrow \text{embedding matrix}$$

$$X \in \mathbb{R}^{T_x \times |V|} \rightarrow \text{one-hot}$$

# Self-Attention

The key idea behing 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.



$$c^{<1>} \quad c^{<2>} \quad c^{<3>} \quad c^{<4>} \quad c^{<5>}$$

$$c^{<3>} = \sum_{t'=1}^{T_x} \alpha^{<3,t'>} v^{<t'>}$$

$$\alpha^{<t,t'>} = softmax(\hat{\alpha}^{<t,t'>})$$

$$\hat{\alpha}^{<t,t'>} = \frac{q^{<t>} \cdot k^{<t'>}}{\sqrt{d_k}}$$

$$q^{<t>} = \mathbf{x}_{epe}^{<t>} W^q$$

$$k^{<t>} = \mathbf{x}_{epe}^{<t>} W^k$$
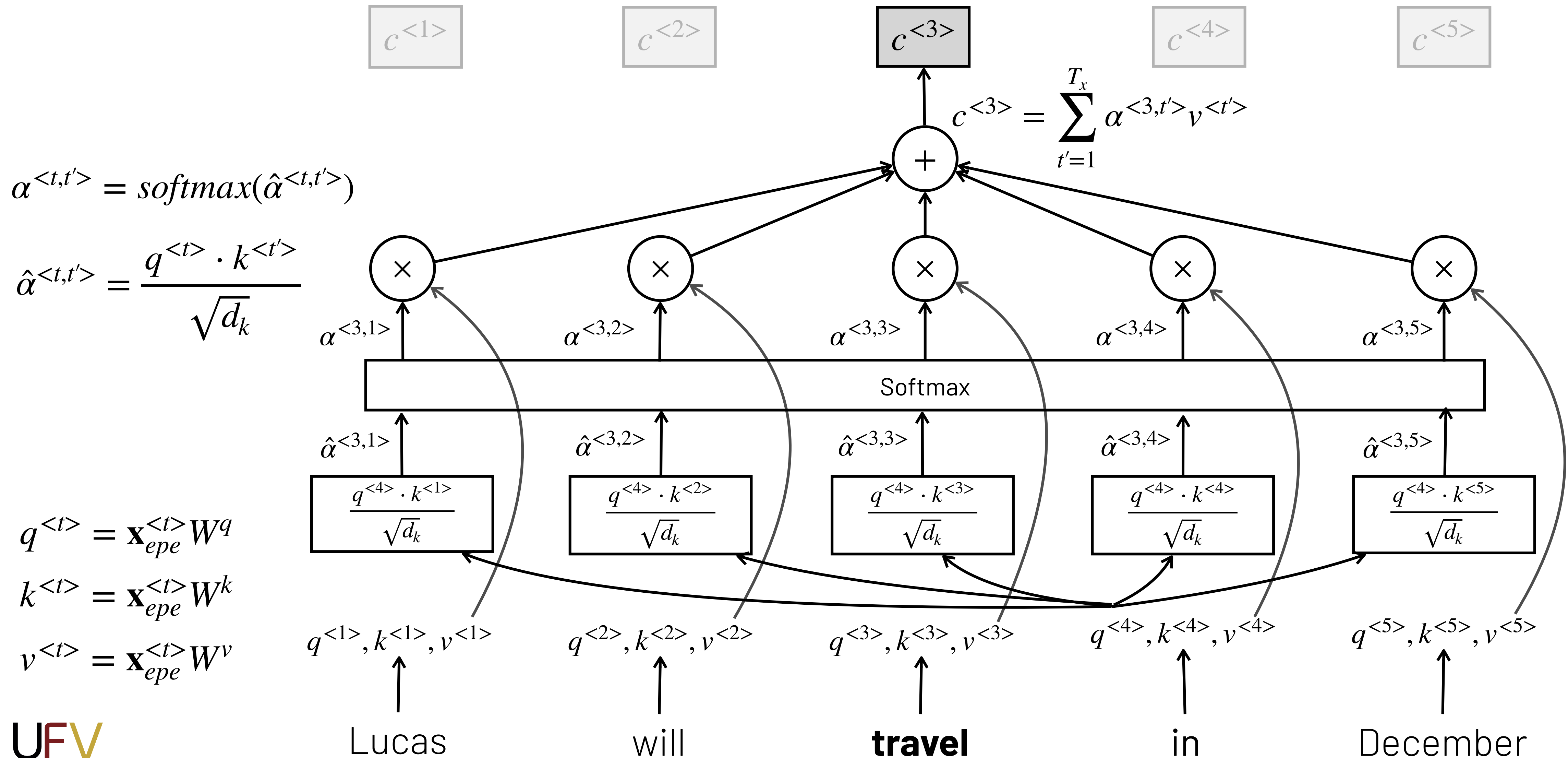
$$v^{<t>} = \mathbf{x}_{epe}^{<t>} W^v$$

$\alpha^{<3,1>} \qquad \alpha^{<3,2>} \qquad \alpha^{<3,3>} \qquad \alpha^{<3,4>} \qquad \alpha^{<3,5>}$

Softmax

$\hat{\alpha}^{<3,1>} \qquad \hat{\alpha}^{<3,2>} \qquad \hat{\alpha}^{<3,3>} \qquad \hat{\alpha}^{<3,4>} \qquad \hat{\alpha}^{<3,5>}$

$$\frac{q^{<4>} \cdot k^{<1>}}{\sqrt{d_k}} \quad \frac{q^{<4>} \cdot k^{<2>}}{\sqrt{d_k}} \quad \frac{q^{<4>} \cdot k^{<3>}}{\sqrt{d_k}} \quad \frac{q^{<4>} \cdot k^{<4>}}{\sqrt{d_k}} \quad \frac{q^{<4>} \cdot k^{<5>}}{\sqrt{d_k}}$$

$q^{<1>}, k^{<1>}, v^{<1>} \quad q^{<2>}, k^{<2>}, v^{<2>} \quad q^{<3>}, k^{<3>}, v^{<3>} \quad q^{<4>}, k^{<4>}, v^{<4>} \quad q^{<5>}, k^{<5>}, v^{<5>}$

Lucas          will          **travel**          in          December

# Self-Attention

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



$C$ $= Attention(Q, K, V) = softmax(\dfrac{QK^T}{\sqrt{(d_k)}})V$

$c^{<1>}$   $c^{<2>}$   $c^{<3>}$   $c^{<4>}$   $c^{<5>}$

$q^{<1>}, k^{<1>}, v^{<1>}$   $q^{<2>}, k^{<2>}, v^{<2>}$   $q^{<3>}, k^{<3>}, v^{<3>}$   $q^{<4>}, k^{<4>}, v^{<4>}$   $q^{<5>}, k^{<5>}, v^{<5>}$

Lucas   will   travel   in   December

UFV

10

# Self-Attention

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

$$Q \times K^T = softmax(\frac{\begin{array}{|c|c|c|} \hline 31 & 27 & 22 \\ \hline 27 & 42 & 31 \\ \hline 22 & 31 & 37 \\ \hline \end{array}}{\sqrt{(d_k)}}) = \begin{array}{c|c|c|c|} & \text{Lucas} & \text{will} & \text{travel} \\ \hline \text{Lucas} & .6 & .2 & .2 \\ \hline \text{will} & .1 & .7 & .2 \\ \hline \text{travel} & .2 & .4 & .4 \\ \hline \end{array} \times V = C$$

$q^{<1>}, q^{<2>}, q^{<3>}$ (rows of $Q$)

$k^{<1>} k^{<2>} k^{<3>}$ (columns of $K^T$)

$q^{<1>}k^{<1>} \quad q^{<1>}k^{<2>} \quad q^{<1>}k^{<3>}$

$v^{<1>}, v^{<2>}, v^{<3>}$ (rows of $V$)

$c^{<1>}, c^{<2>}, c^{<3>}$ (rows of $C$)

Attention weights of word $i$ to word $j$ (normalized by line)

$$Q = X_{epe} \times W^q \in \mathbb{R}^{d_e \times d_q}$$

$$K = X_{epe} \times W^k \in \mathbb{R}^{d_e \times d_k}$$

$$V = X_{epe} \times W^v \in \mathbb{R}^{d_e \times d_v}$$

$$X_{epe} \in \mathbb{R}^{T_x \times d_e}$$

Lucas, will, travel

$$d_e = d_q = d_k = d_v = 3$$

The attention layer receives the embedded sequence $X_{epe}$ as input

The sizes of key and query have to be the same. But embeddings and value typically also have same sizes.

UFV

11

# Multi-Head Attention

$$Multihead(Q, K, V) = Concat(C_1, C_2, \ldots, C_h)W^o$$

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

$$[]$$

$$W_i^q, W_i^k, W_i^v \quad C_i \quad = Attention(Q, K, V) = softmax(\frac{QK^T}{\sqrt{(d_k)}})V$$

$c_i^{<1>}$     $c_i^{<2>}$     $c_i^{<3>}$     $c_i^{<4>}$     $c_i^{<5>}$

$q^{<1>}, k^{<1>}, v^{<1>}$   $q^{<2>}, k^{<2>}, v^{<2>}$   $q^{<3>}, k^{<3>}, v^{<3>}$   $q^{<4>}, k^{<4>}, v^{<4>}$   $q^{<5>}, k^{<5>}, v^{<5>}$

$\mathbf{x}_{epe}^{<1>}$     $\mathbf{x}_{epe}^{<2>}$     $\mathbf{x}_{epe}^{<3>}$     $\mathbf{x}_{epe}^{<4>}$     $\mathbf{x}_{epe}^{<5>}$

Lucas      will      travel      in      December

UFV

# Multi-Head Attention

$$Multihead(Q, K, V) = Concat(C_1, C_2, C_3)W^o$$



$C_1$ $C_2$ $C_3$

$W^o \in \mathbb{R}^{d_{model} \times d_{model}}$

$\times$ $= C \in \mathbb{R}^{T_x \times d_{model}}$

$C_1 = softmax(\frac{Q_1 K_1^T}{\sqrt{(d_k)}})V_1$

$C_2 = softmax(\frac{Q_2 K_2^T}{\sqrt{(d_k)}})V_2$

$C_3 = softmax(\frac{Q_3 K_3^T}{\sqrt{(d_k)}})V_3$

$Q_1 = X_{epe}W_1^q$　　$K_1 = X_{epe}W_1^k$　　$V_1 = X_{epe}W_1^v$　　$Q_2 = X_{epe}W_2^q$　　$K_2 = X_{epe}W_2^k$　　$V_2 = X_{epe}W_2^v$　　$Q_3 = X_{epe}W_3^q$　　$K_3 = X_{epe}W_3^k$　　$V_3 = X_{epe}W_3^v$

Number of heads

$h = 3$

$d_k = \dfrac{d_{model}}{h} = 3$

$W^q \in \mathbb{R}^{d_{model} \times d_{model}}$　　$W^k \in \mathbb{R}^{d_{model} \times d_{model}}$　　$W^v \in \mathbb{R}^{d_{model} \times d_{model}}$　　$W^o \in \mathbb{R}^{d_{model} \times d_{model}}$

$W_1^q$　$W_2^q$　$W_3^q$　　　$W_1^k$　$W_2^k$　$W_3^k$　　　$W_1^v$　$W_2^v$　$W_v^k$

$T_x = 3$
(sequence length)

$d_{model} = 9$

Lucas

will

travel

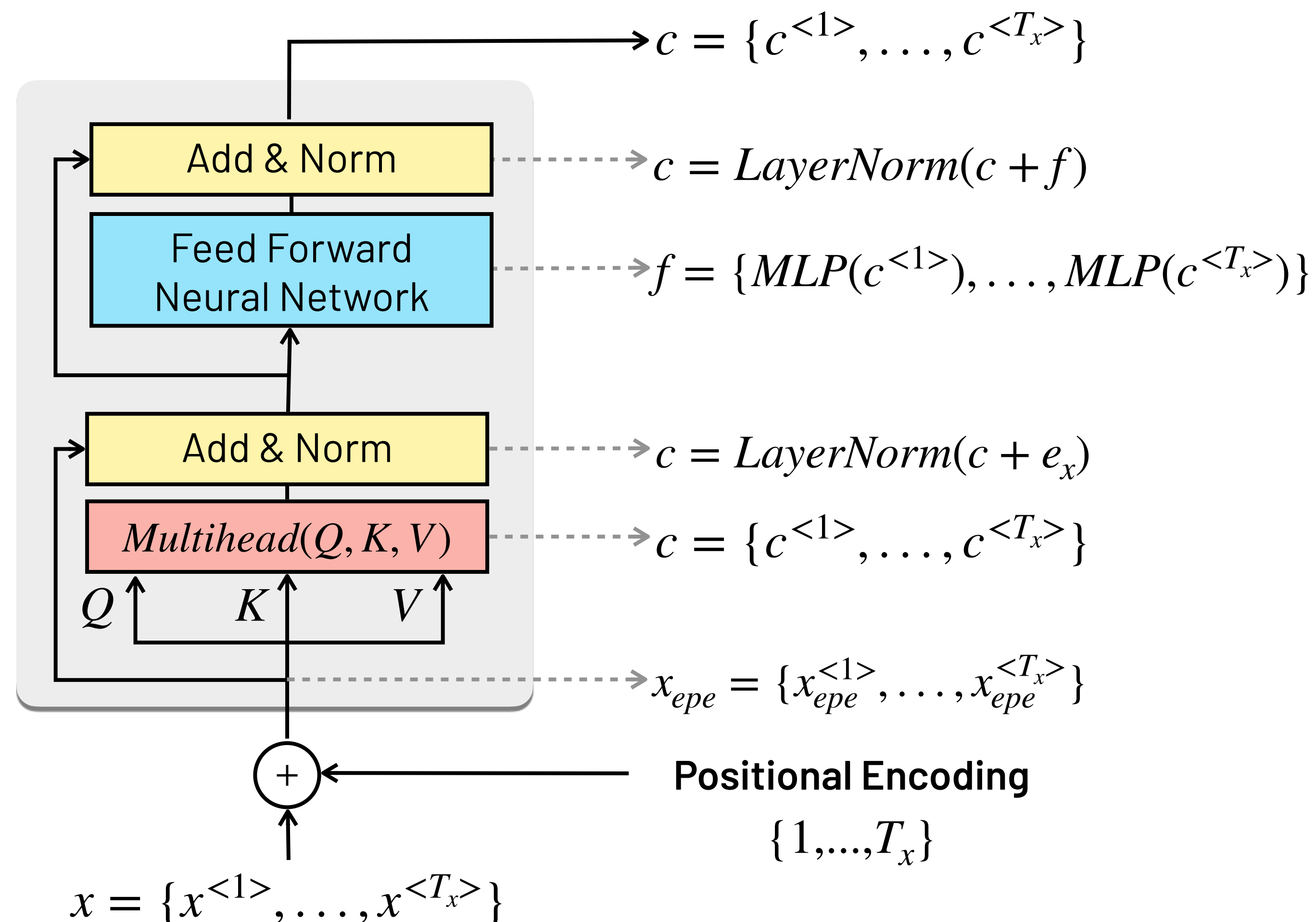$X_{epe} \in \mathbb{R}^{T_x \times d_{model}}$

# Encoder

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

**Ouput**: a contextual representation $C = \{c^{<1>}, \ldots, c^{<T_x>}\}$ of $x$

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

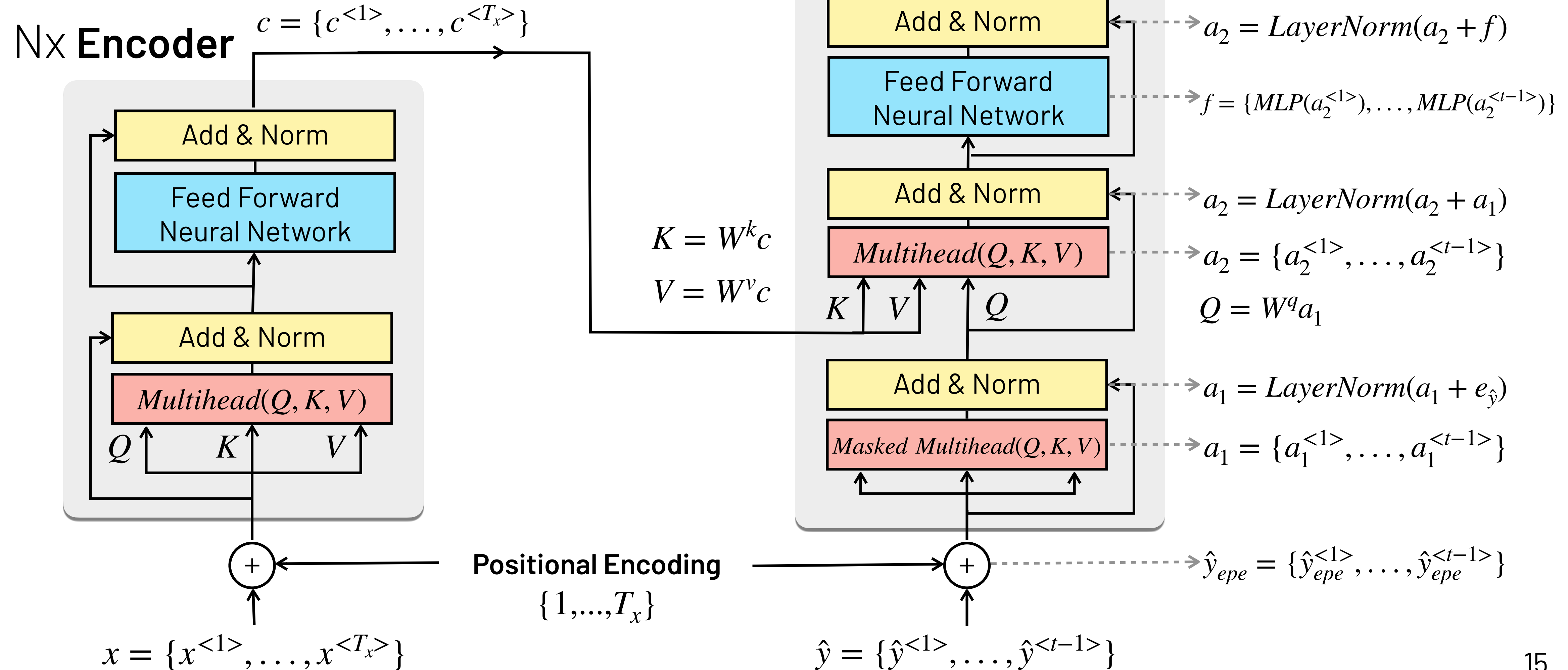Both are normalized with Layer Norm (**Norm**) and conneced with a residual connection (**Add**)

$c = \{c^{<1>}, \ldots, c^{<T_x>}\}$

| Add & Norm | $\quad\dashrightarrow c = LayerNorm(c + f)$ |

Feed Forward Neural Network $\quad\dashrightarrow f = \{MLP(c^{<1>}), \ldots, MLP(c^{<T_x>})\}$

Add & Norm $\quad\dashrightarrow c = LayerNorm(c + e_x)$

$Multihead(Q, K, V) \quad\dashrightarrow c = \{c^{<1>}, \ldots, c^{<T_x>}\}$

$Q \quad K \quad V$

$x_{epe} = \{x_{epe}^{<1>}, \ldots, x_{epe}^{<T_x>}\}$

$+$

**Positional Encoding**
$\{1, \ldots, T_x\}$

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

UFV

# Decoder

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

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

$\hat{y}^{<t>}$

Softmax

fc

**Decoder** Nx

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

Nx **Encoder**

$c = \{c^{<1>}, \ldots, c^{<T_x>}\}$

Add & Norm

$a_2 = LayerNorm(a_2 + f)$

Feed Forward
Neural Network

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

Add & Norm

$a_2 = LayerNorm(a_2 + a_1)$

$K = W^k c$

$V = W^v c$

$Multihead(Q, K, V)$

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

$K$  $V$    $Q$

$Q = W^q a_1$

Add & Norm

Feed Forward
Neural Network

Add & Norm

$Multihead(Q, K, V)$

$Q$    $K$    $V$

Add & Norm

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

$Masked\ Multihead(Q, K, V)$

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

**Positional Encoding**
$\{1, \ldots, T_x\}$

$+$

$\hat{y}_{epe} = \{\hat{y}_{epe}^{<1>}, \ldots, \hat{y}_{epe}^{<t-1>}\}$
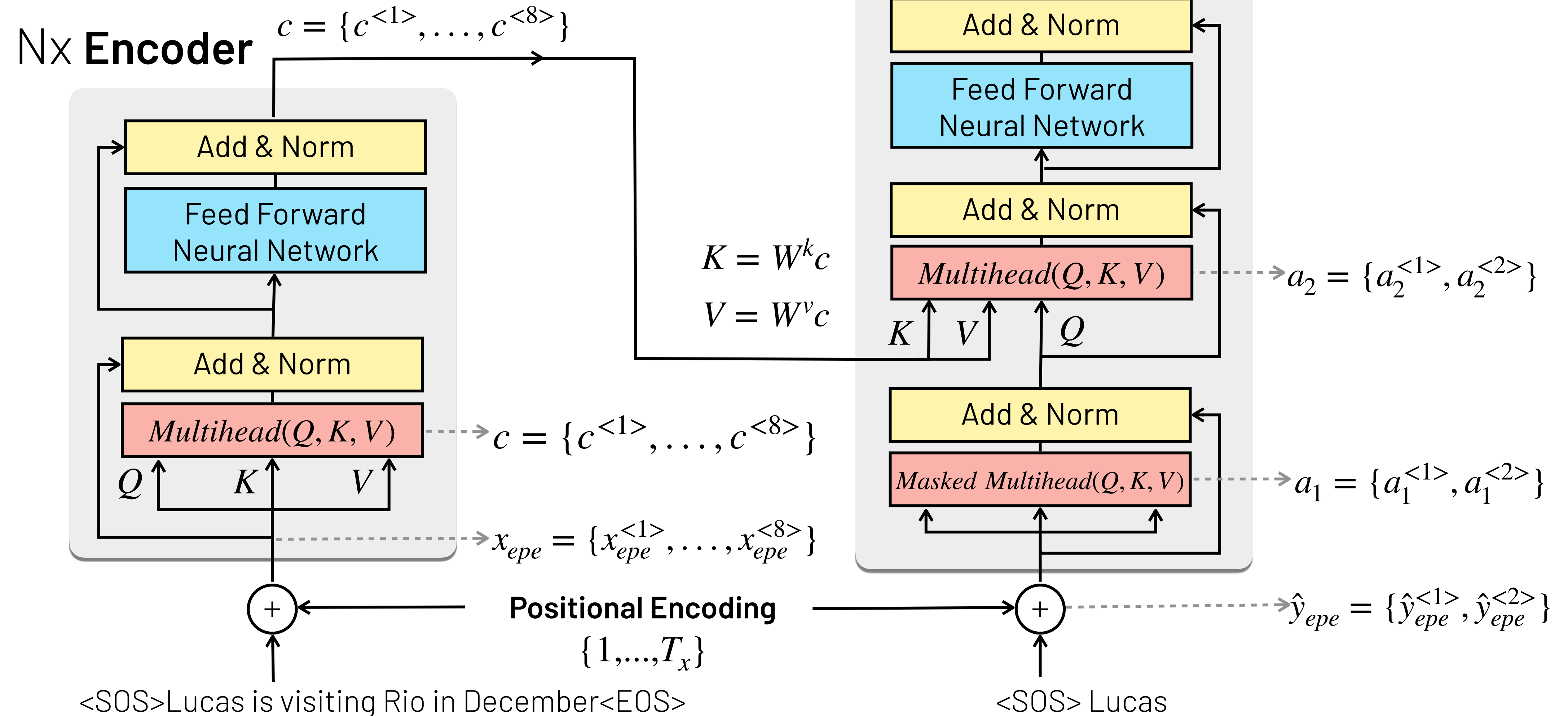
$x = \{x^{<1>}, \ldots, x^{<T_x>}\}$

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

# Inference Example

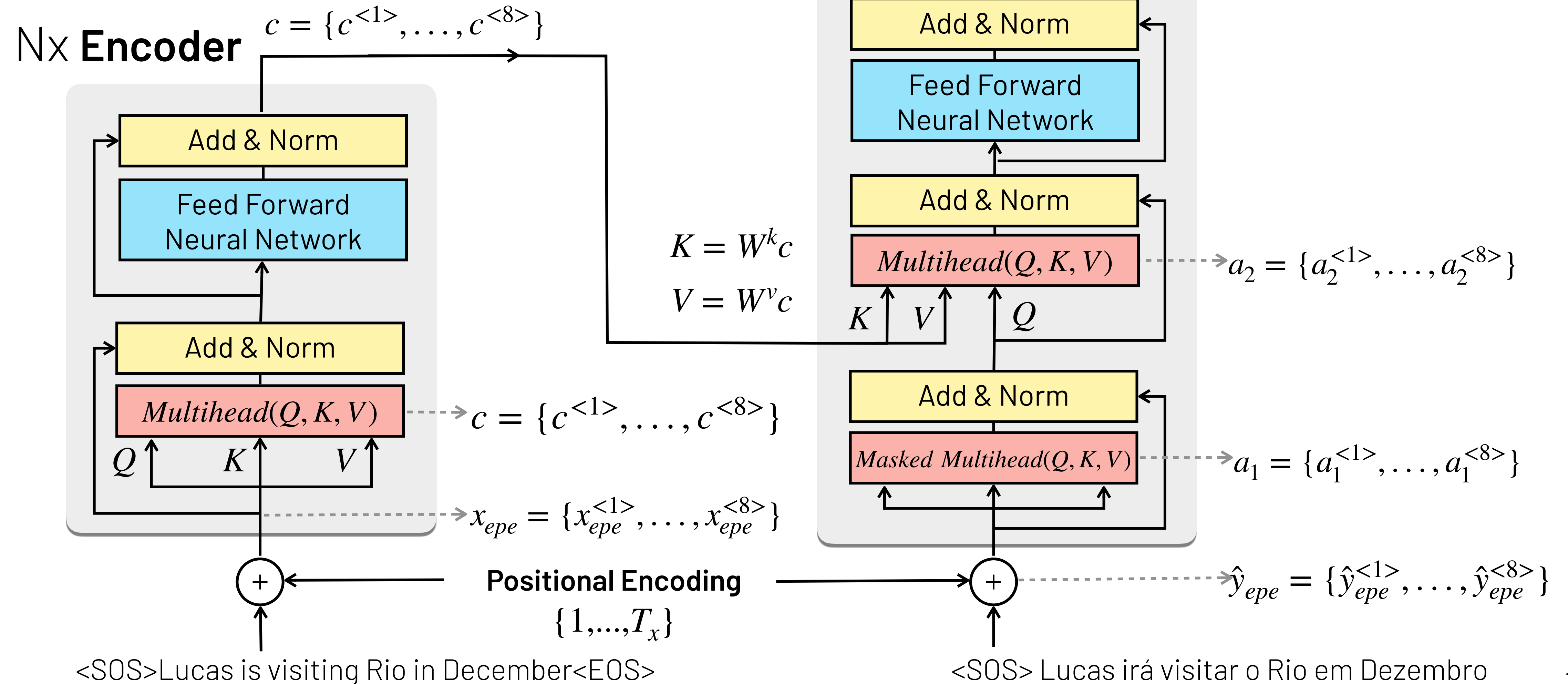At inference time, the model is autoregressive, i.e., it generates one word at a time.

$$\hat{y} = \text{Lucas}$$



**Decoder** Nx

Nx **Encoder**

$$c = \{c^{<1>}, \ldots, c^{<8>}\}$$

Softmax

fc

$$a_2 = \{a_2^{<1>}\}$$

Add & Norm

Feed Forward Neural Network

Add & Norm

Add & Norm

Feed Forward Neural Network

$$K = W^k c$$
$$V = W^v c$$

Add & Norm

$Multihead(Q, K, V)$

$$a_2 = \{a_2^{<1>}\}$$

$K$  $V$   $Q$

Add & Norm

$Multihead(Q, K, V)$

$$c = \{c^{<1>}, \ldots, c^{<8>}\}$$

$Q$   $K$   $V$

$Masked\ Multihead(Q, K, V)$

$$a_1 = \{a_1^{<1>}\}$$

$$x_{epe} = \{x_{epe}^{<1>}, \ldots, x_{epe}^{<8>}\}$$

$+$

$+$

$$\hat{y}_{epe} = \{\hat{y}_{epe}^{<1>}\}$$

**Positional Encoding**
$$\{1, \ldots, T_x\}$$

<SOS>Lucas is visiting Rio in December<EOS>

<SOS>

UFV

16

# Inference Example

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

$\hat{y}$ = Lucas irá

**Softmax**

**fc**

**Decoder** Nx

$a_2 = \{a_2^{<1>}, a_2^{<2>}\}$

$c = \{c^{<1>}, \dots, c^{<8>}\}$

Nx **Encoder**

**Add & Norm**

**Feed Forward Neural Network**

**Add & Norm**

**Multihead**(Q, K, V)

$Q$   $K$   $V$

$c = \{c^{<1>}, \dots, c^{<8>}\}$

$K = W^k c$

$V = W^v c$

**Add & Norm**

**Feed Forward Neural Network**

**Add & Norm**

*Multihead*(Q, K, V)

$K$   $V$        $Q$

$a_2 = \{a_2^{<1>}, a_2^{<2>}\}$

**Add & Norm**

*Masked Multihead*(Q, K, V)

$a_1 = \{a_1^{<1>}, a_1^{<2>}\}$

$x_{epe} = \{x_{epe}^{<1>}, \dots, x_{epe}^{<8>}\}$

$+$

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

$+$

$\hat{y}_{epe} = \{\hat{y}_{epe}^{<1>}, \hat{y}_{epe}^{<2>}\}$

<SOS>Lucas is visiting Rio in December<EOS>

<SOS> Lucas

UFV

16

# Inference Example

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

$\hat{y}$ = Lucas irá visitar o Rio em Dezembro<EOS>

**Softmax**

**Decoder** Nx

**fc**

$a_2 = \{a_2^{<1>}, \ldots, a_2^{<8>}\}$

Nx **Encoder**

$c = \{c^{<1>}, \ldots, c^{<8>}\}$

**Add & Norm**

**Add & Norm**

**Feed Forward Neural Network**

**Feed Forward Neural Network**

**Add & Norm**

$K = W^k c$

**Add & Norm**

$V = W^v c$

$Multihead(Q, K, V)$

$a_2 = \{a_2^{<1>}, \ldots, a_2^{<8>}\}$

**Add & Norm**

$K$ $V$ $Q$

$Multihead(Q, K, V)$

$c = \{c^{<1>}, \ldots, c^{<8>}\}$

**Add & Norm**

$Q$ $K$ $V$

$Masked\ Multihead(Q, K, V)$

$a_1 = \{a_1^{<1>}, \ldots, a_1^{<8>}\}$

$x_{epe} = \{x_{epe}^{<1>}, \ldots, x_{epe}^{<8>}\}$

+

**Positional Encoding** $\{1, \ldots, T_x\}$

+

$\hat{y}_{epe} = \{\hat{y}_{epe}^{<1>}, \ldots, \hat{y}_{epe}^{<8>}\}$

<SOS>Lucas is visiting Rio in December<EOS>

<SOS> Lucas irá visitar o Rio em Dezembro

UFV

16

# Positional Encoding

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

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

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

$$x_{epe} = \{x_{epe}^{<1>}, \ldots, x_{epe}^{<T_x>}\}$$

**Positional Encoding** $=$

$$PE_{(t,2i)} = sin(\frac{t}{10000^{\frac{2i}{d}}})$$

$$PE_{(t,2i+1)} = cos(\frac{t}{10000^{\frac{2i}{d}}})$$

$$E \quad x_e = \{x_e^{<1>}, \ldots, x_e^{<T_x>}\}$$

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

$t = 1$

$i = 0$
$i = 1$
$i = 2$
$i = 3$

$= \quad + \quad d_{model} = 4$

$x_{epe}^{<1>} \quad x_e^{<1>} \quad pe^{<1>}$

# Positional Encoding

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

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



$sin$
$i = 0$

$cos$
$i = 1$

$sin$
$i = 2$

$cos$
$i = 3$

$d = 4$

$t:$   $x_e^{<1>}$   $pe^{<1>}$     $x_e^{<2>}$   $pe^{<2>}$     $x_e^{<3>}$   $pe^{<3>}$     $x_e^{<4>}$   $pe^{<4>}$     $x_e^{<5>}$   $pe^{<5>}$     $x_e^{<6>}$   $pe^{<6>}$

Lucas      is      visiting      **Rio**      in      December

UFV

# Training with Masked Multi-head Attention

$$\hat{y}^{<1>} \quad \hat{y}^{<2>} \quad \hat{y}^{<3>} \quad \hat{y}^{<4>} \quad \hat{y}^{<5>}$$

$\hat{y} =$ Lucas    irá    visitar    o    ___

**Encoder**

**Decoder**

*Masked Multihead*

- During training, the model **does not** produce $\hat{y}$ one token at at time;

- Instead, the output $\hat{y}$ is produced once, in parallel with vectorization.

- Thus, the full target sentence $y = \{y^{<1>}, \ldots, y^{<T_y>}\}$ is given as input to the Decoder, which uses a **mask** to predict $\hat{y}^{<t>}$, so it does not attent to the future tokens ($\hat{y}^{<t'>}$, $t' > t$). For example, for $t = 5$:

$x =$ <SOS>Lucas is visiting Rio in December <EOS>

$y =$ <SOS>   Lucas    irá    visitar    o    Rio    em   Dezembro

$$y^{<1>} \quad y^{<2>} \quad y^{<3>} \quad y^{<4>} \quad y^{<5>} \quad y^{<6>} \quad y^{<7>} \quad y^{<8>}$$

UFV

# Training with Masked Multi-head Attention

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

$Q$  $K^T$

| 31 | 27 | 22 |
| 27 | 42 | 31 |
| 22 | 31 | 37 |

$\times$ = $\overline{\sqrt{(d_k)}}$

**Mask**

| 0 | -inf | -inf |
| 0 | 0 | -inf |
| 0 | 0 | 0 |

$+$

$= softmax($

| 31 | -inf | -inf |
| 27 | 42 | -inf |
| 22 | 31 | 37 |

$)$ $= \overline{\sqrt{(d_k)}}$

| .1 | .0 | .0 |
| .3 | .7 | .0 |
| .3 | .3 | .4 |

$Q$  $X_{epe}$  $W^q \in \mathbb{R}^{d_e \times d_q}$

= $\times$

$K$  $X_{epe}$  $W^k \in \mathbb{R}^{d_e \times d_k}$

= $\times$

$V$  $X_{epe}$  $W^v \in \mathbb{R}^{d_e \times d_v}$

= $\times$

The attention layer receives the embedded sequence $X_{epe}$ as input

$X_{epe} \in \mathbb{R}^{T_x \times d_e}$

Lucas

will

travel

$d_e = d_q = d_k = d_v = 3$

The sizes of key and query have to be the same. But embeddings and value typically also have same sizes.

UFV

# Next Lecture

**L17**: Transformers (Part II)

Case studies of transformers: BERT and GPT

UFV