



Deep Learning

Transformers

Dr. Mehran Safayani

safayani@iut.ac.ir

safayani.iut.ac.ir



<https://www.aparat.com/mehran.safayani>



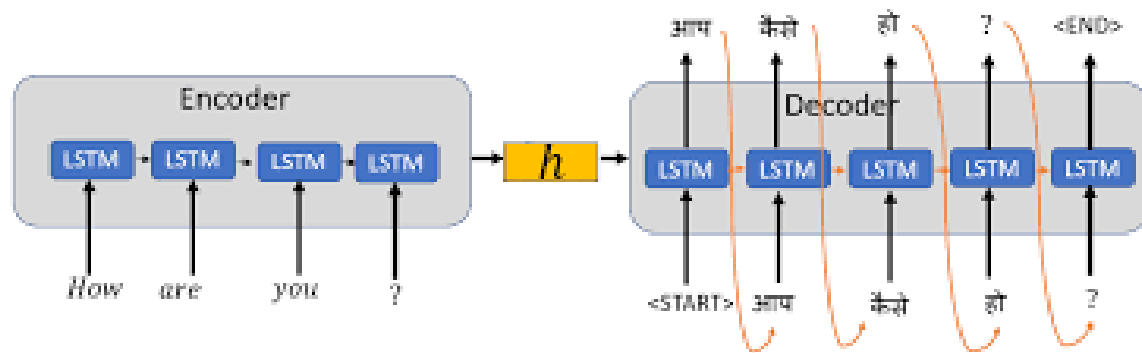
https://github.com/safayani/deep_learning_course



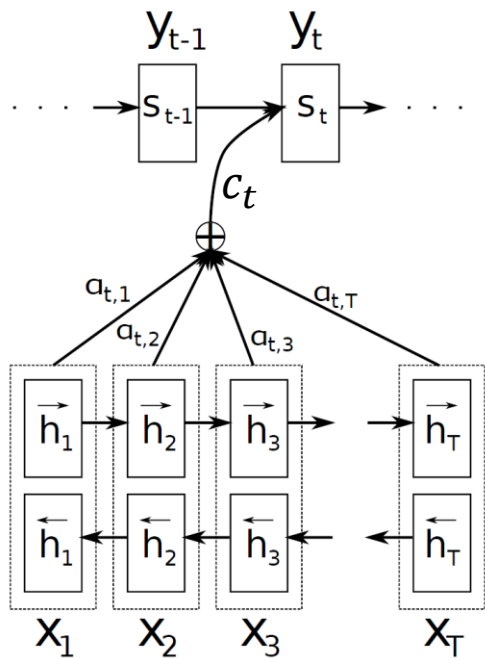
Attention model

A potential issue with this encoder–decoder approach is that a neural network needs to be able to compress all the necessary information of a source sentence into a fixed-length vector.

This may make it difficult for the neural network to cope with long sentences, especially those that are longer than the sentences in the training corpus.



Attention model



$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j.$$

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})},$$

$$e_{ij} = a(s_{i-1}, h_j)$$

a is a feed forward NN

$$s_i = f(s_{i-1}, y_{i-1}, c_i).$$

$$p(y_i | y_1, \dots, y_{i-1}, \mathbf{x}) = g(y_{i-1}, s_i, c_i),$$

Attention model summary

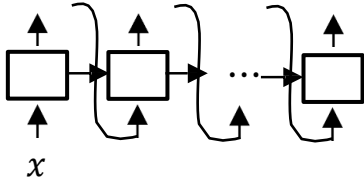
Each time the proposed model generates a word in a translation, it (soft-)searches for a set of positions in a source sentence where the most relevant information is concentrated. The model then predicts a target word based on the context vectors associated with these source positions and all the previous generated target words.

The most important distinguishing feature of this approach from the basic encoder–decoder is that it does not attempt to encode a whole input sentence into a single fixed-length vector. Instead, it encodes the input sentence into a sequence of vectors and chooses a subset of these vectors adaptively while decoding the translation.

Transformers Motivation

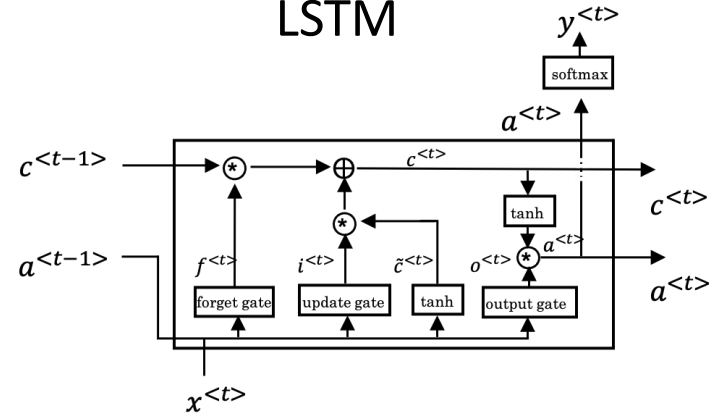
Increased complexity,
sequential

RNN



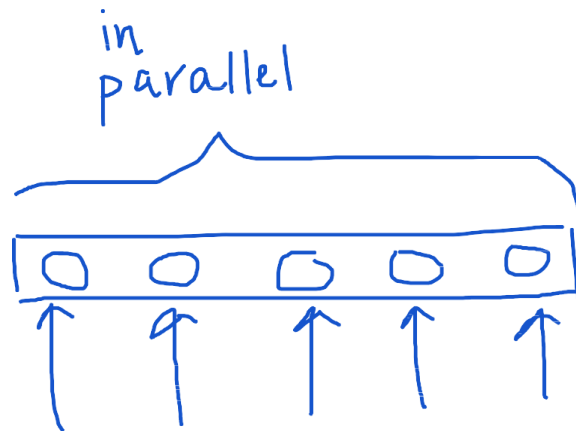
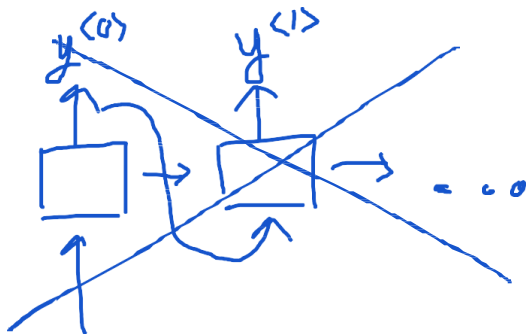
GRU

LSTM



Transformers Intuition

- Attention + CNN
 - Self-Attention
 - Multi-Head Attention



Self-Attention Intuition

$A(q, K, V)$ = attention-based vector representation of a word

RNN Attention

$$\alpha^{<t, t'>} = \frac{\exp(<t, t'>)}{\sum_{t'=1}^T \exp(<t, t'>)}$$

$$q^t = W^Q x^t$$

$$k^t = W^K x^t$$

Transformers Attention

$$A(q, K, V) = \sum_i \frac{\exp(<q \cdot k^{<i>}>)}{\sum_j \exp(<q \cdot k^{<j>}>)} v^{<i>}$$

$$v^t = W^V x^t$$

Query, Key, Value

The key/value/query concept is analogous to retrieval systems. For example, when you search for videos on internet, the search engine will map your **query** (text in the search bar) against a set of **keys** (video title, description, etc.) associated with candidate videos in their database, then present you the best matched videos (**values**).

Self-Attention Intuition

$A(q, K, V)$ = attention-based vector representation of a word

RNN Attention

$$\alpha^{<t, t'>} = \frac{\exp(<t, t'>)}{\sum_{t'=1}^T \exp(<t, t'>)}$$

Transformers Attention

$$A(q, K, V) = \sum_i \frac{\exp(<q \cdot k^{<i>}>)}{\sum_j \exp(<q \cdot k^{<j>}>)} v^{<i>}$$

$$q^t = W^Q x^t$$

$$k^t = W^K x^t$$

$$v^t = W^V x^t$$

q^3 intuition: what's happening in Africa?

$x^{<1>}$
Jane

$x^{<2>}$
visite

$x^{<3>}$
l'Afrique

$x^{<4>}$
en

$x^{<5>}$
septembre

Self-Attention Intuition

$$q^t = W^Q x^t$$

$$k^t = W^K x^t$$

$$v^t = W^V x^t$$

q^3 : what's happening in Africa?

$k^{<1>}$: person

$v^{<1>}$: jane embedding

$k^{<2>}$: action

$v^{<2>}$: visit embedding

Query (Q)	Key (K)	Value (V)
$q^{<1>}$	$k^{<1>}$	$v^{<1>}$
$q^{<2>}$	$k^{<2>}$	$v^{<2>}$
$q^{<3>}$	$k^{<3>}$	$v^{<3>}$
$q^{<4>}$	$k^{<4>}$	$v^{<4>}$
$q^{<5>}$	$k^{<5>}$	$v^{<5>}$

$x^{<1>}$
Jane

$x^{<2>}$
visite

$x^{<3>}$
l'Afrique

$x^{<4>}$
en

$x^{<5>}$
septembre

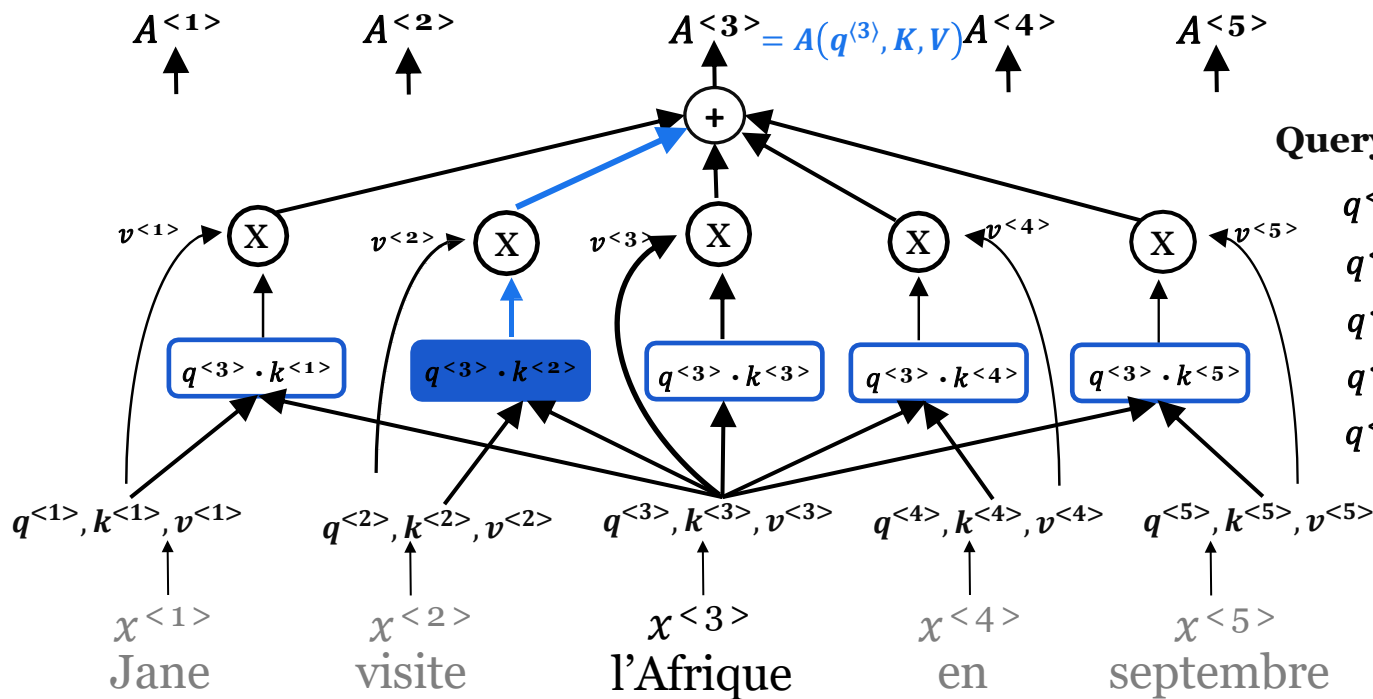
Self-Attention

If you put all of these computation together:

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

$$A(q, K, V) = \sum_i \frac{\exp(q \cdot k^{<i>})}{\sum_j \exp(q \cdot k^{<j>})} v^{<i>}$$

Contextual Embedding



$$q^{<3>} = W^Q \cdot x^{<3>}$$

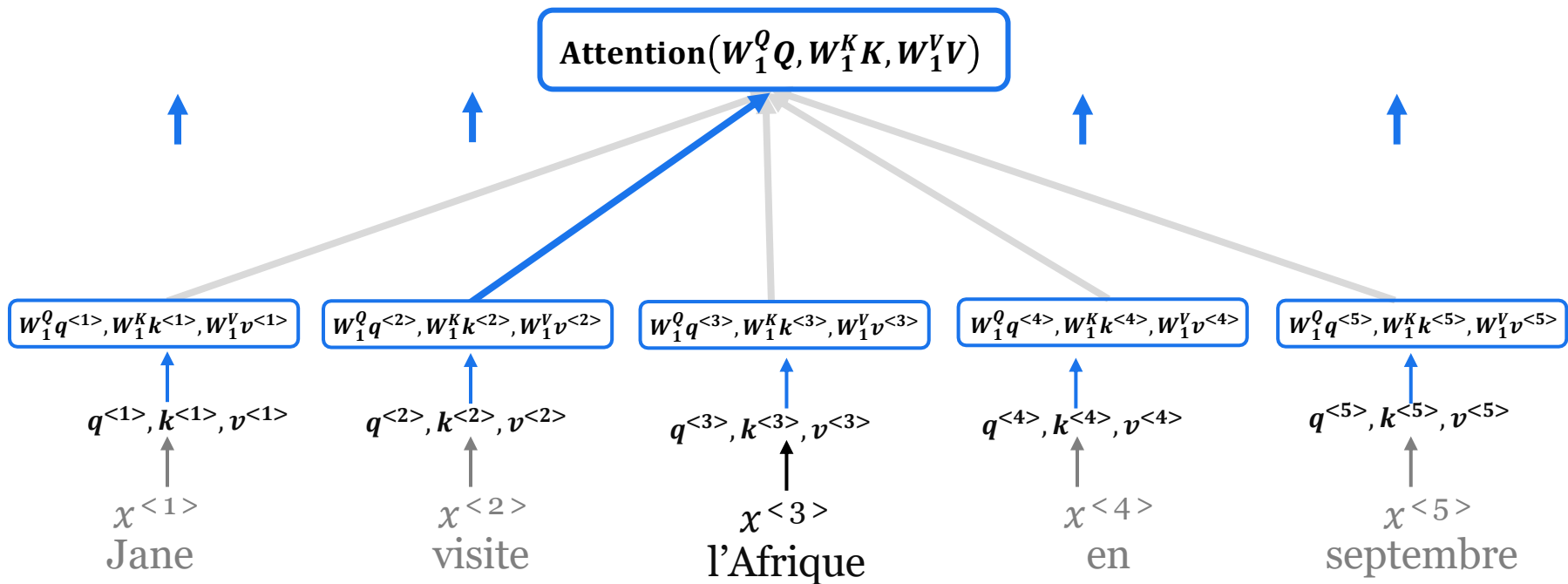
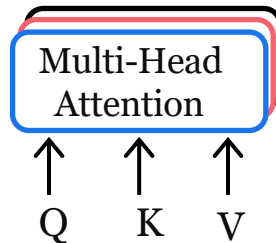
$$k^{<3>} = W^K \cdot x^{<3>}$$

$$v^{<3>} = W^V \cdot x^{<3>}$$

Multi-Head Attention

first head: w_1^Q, w_1^K, w_1^V - What's happening in Africa?

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



Multi-Head Attention

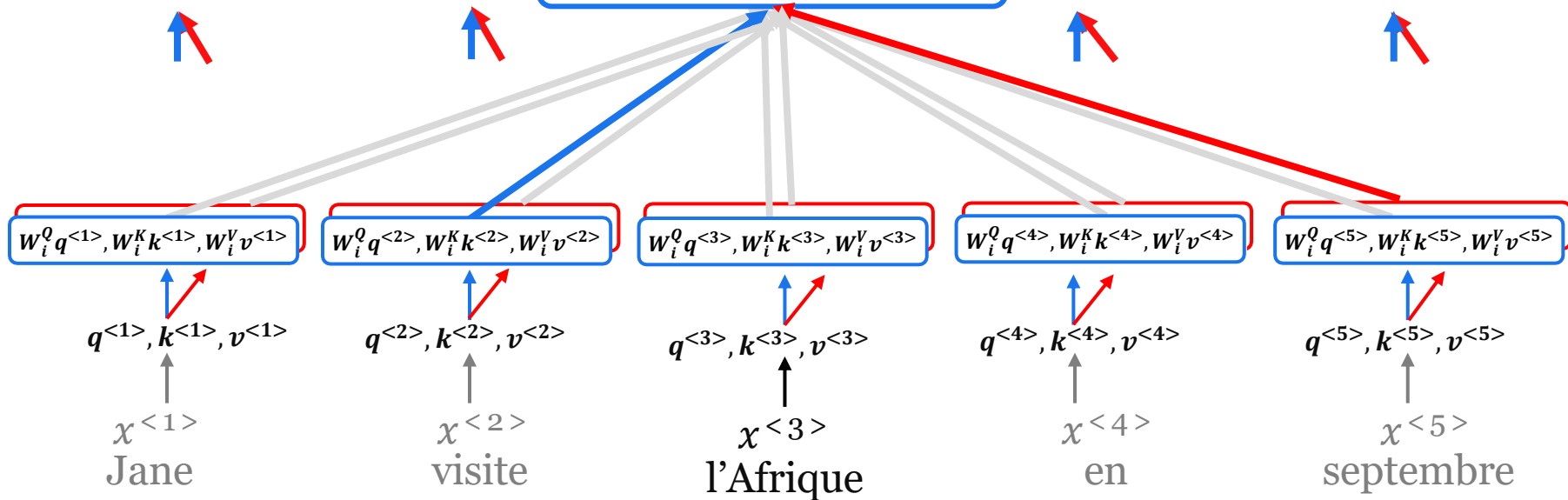
second head: W_2^Q, W_2^K, W_2^V - When is sth happening in Africa?

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

$$\text{Attention}(W_i^Q Q, W_i^K K, W_i^V V)$$

Multi-Head
Attention

↑ Q ↑ K ↑ V



Multi-Head Attention

third head: W_3^Q, W_3^K, W_3^V - Who has sth to do in Africa?

$$\text{MultiHead}(Q, K, V) = \text{concat}(\text{head}_1 \text{ head}_2 \dots \text{head}_h) W_o$$

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

$$\text{head}_i = \text{Attention}(W_i^Q Q, W_i^K K, W_i^V V)$$

Multi-Head
Attention

Q
K
V

$h = \# \text{heads}$

$$\text{Attention}(W_i^Q Q, W_i^K K, W_i^V V)$$

$q^{<1>}, k^{<1>}, v^{<1>}$

$x^{<1>}$
Jane

$q^{<2>}, k^{<2>}, v^{<2>}$

$x^{<2>}$
visite

$q^{<3>}, k^{<3>}, v^{<3>}$

$x^{<3>}$
l'Afrique

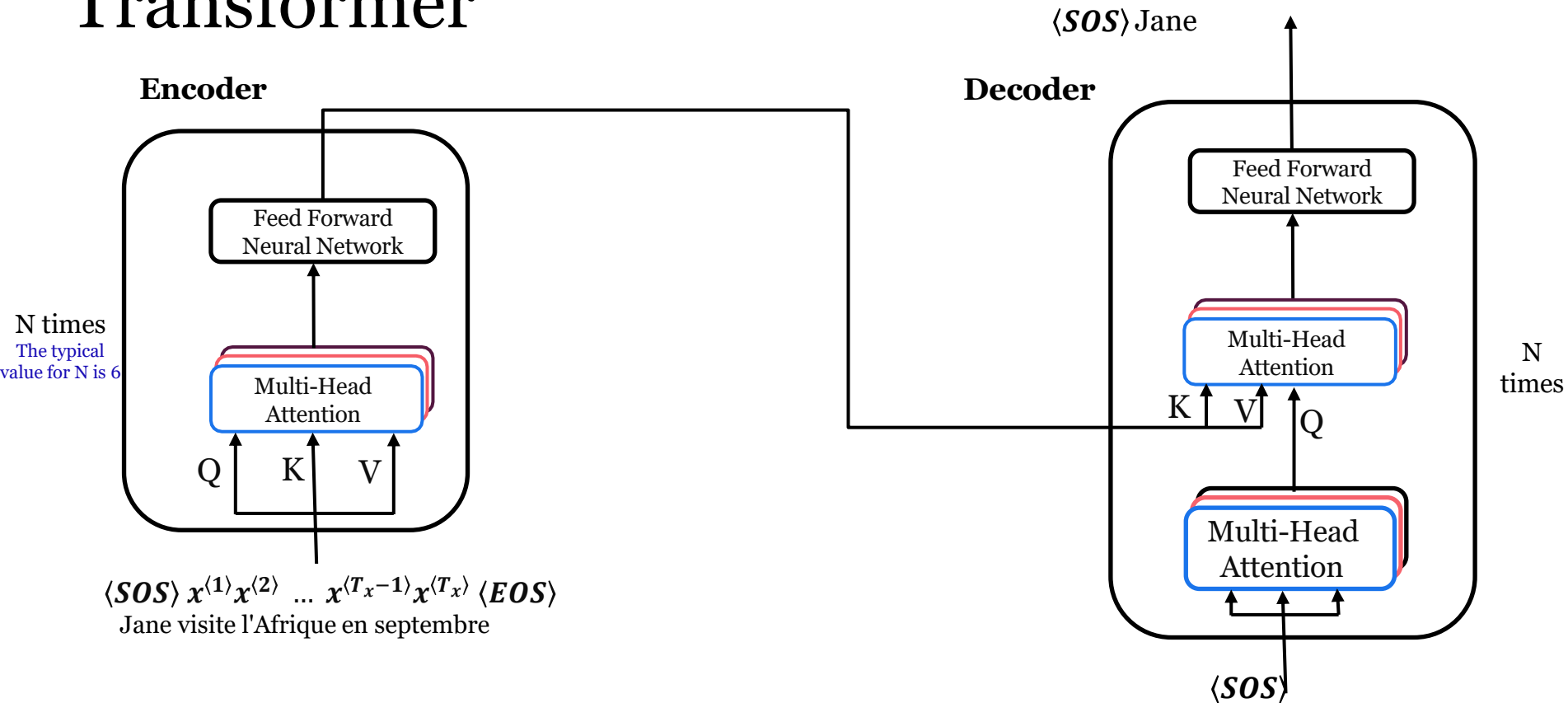
$q^{<4>}, k^{<4>}, v^{<4>}$

$x^{<4>}$
en

$q^{<5>}, k^{<5>}, v^{<5>}$

$x^{<5>}$
septembre

Transformer

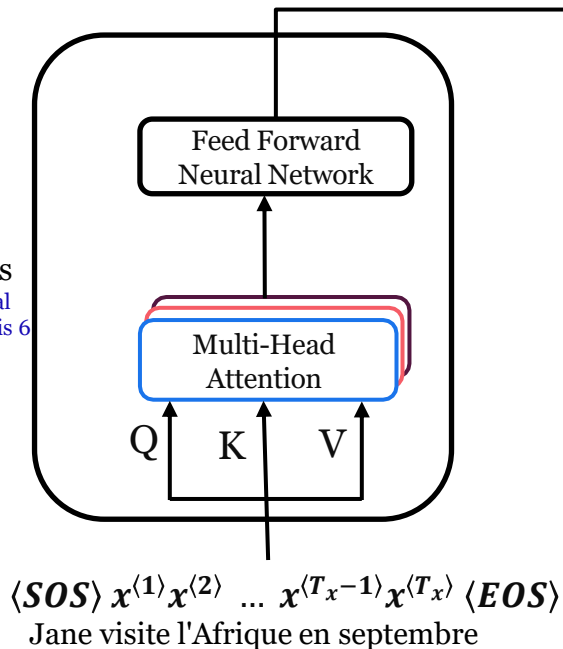


- The translation will start with a start of sentence token ($\langle \text{SOS} \rangle$) and so the start of sentence token gets fed into this multi-head attention block.
- Just this one token, $\langle \text{SOS} \rangle$ is used to compute Q , K and V for this multi-head attention block.

Transformer

Encoder

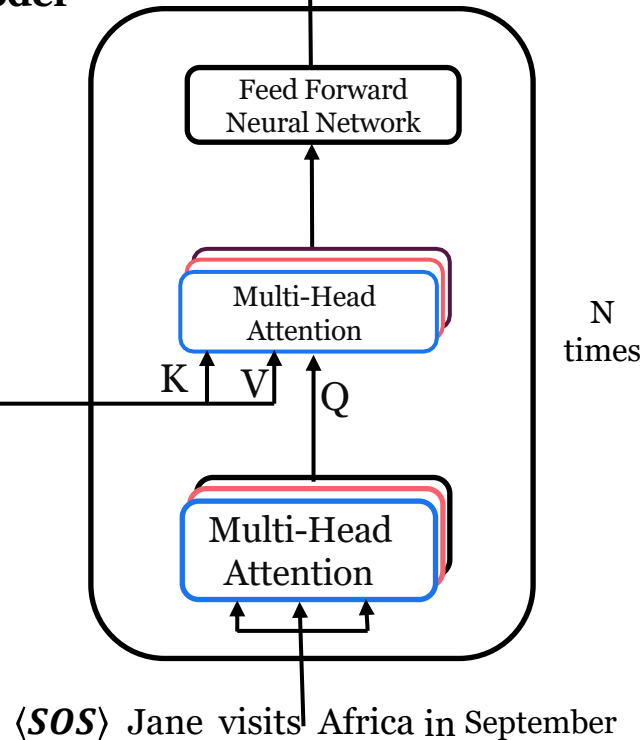
N times
The typical
value for N is 6



$\langle SOS \rangle$ Jane visits Africa in September $\langle EOS \rangle$

Decoder

N
times



Transformer Details: Positional Embedding

- Positional encoding is necessary to convey word positions in self-attention.
- Self-attention equations don't inherently indicate word positions.
- Position within a sentence holds significant importance.
- Sine and cosine equations are used for positional encoding.
- Encoding position helps capture meaningful context.
- Preserving positional information enhances sentence understanding.

Positional Encoding

- First idea: add a value in range $[0,1]$ to the time steps in which 0 means the first word and 1 is the last time-step

Drawbacks:

- you can't figure out how many words are present within a specific range.
- time-step delta doesn't have consistent meaning across different sentences.

- Second idea: assign a number to each time-step linearly. That is, the first word is given "1", the second word is given "2", and so on.

Drawbacks:

- The values could get quite large
- our model can face sentences longer than the ones in training.
- our model may not see any sample with one specific length which would hurt generalization of our model.

Positional Encoding

Ideally, the following criteria should be satisfied:

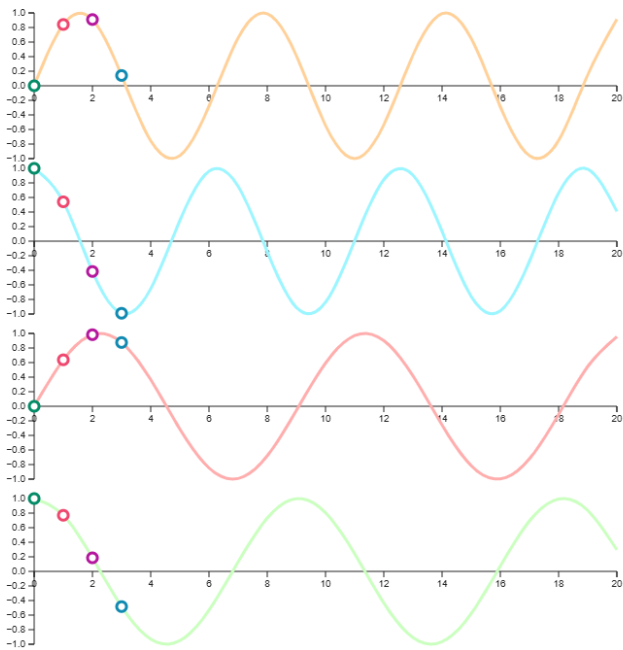
- It should output a unique encoding for each time-step (word's position in a sentence)
- Distance between any two time-steps should be consistent across sentences with different lengths.
- Our model should generalize to longer sentences without any efforts. Its values should be bounded.
- It must be deterministic

Positional Encoding

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

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

Positional encoding visualization



p0	p1	p2	p3
0.000	0.841	0.909	0.141
1.000	0.540	-0.416	-0.990
0.000	0.638	0.983	0.875
1.000	0.770	0.186	-0.484

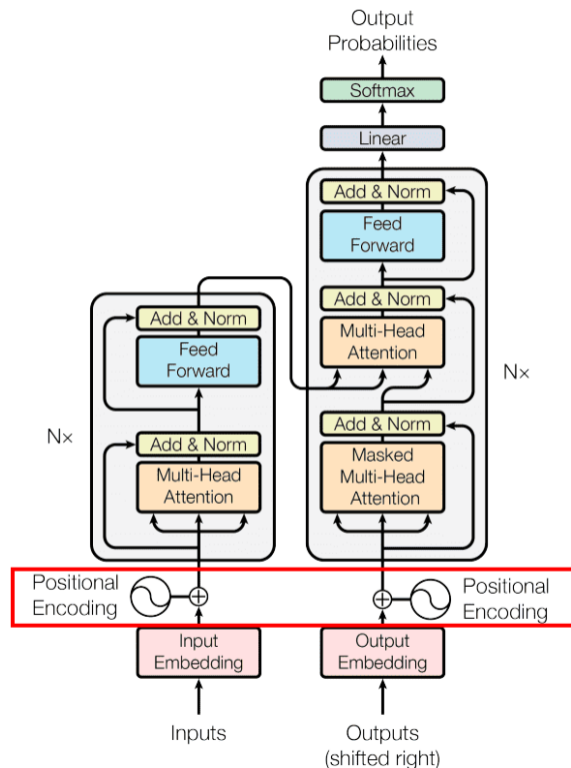
Positional Encoding

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

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

Settings: d = 50

The value of each positional encoding depends on the *position* (*pos*) and *dimension* (*d*). We calculate result for every *index* (*i*) to get the whole vector.



Layer Normalizaion

batch normalization drawbacks:

- batch normalization requires running averages of the summed input statistics. In feed-forward networks with fixed depth, it is straightforward to store the statistics separately for each hidden layer. However, the summed inputs to the recurrent neurons in a recurrent neural network (RNN) often vary with the length of the sequence so applying batch normalization to RNNs appears to require different statistics for different time-steps.
- BN cannot be applied to online learning tasks or to extremely large distributed models where the mini-batches have to be small.

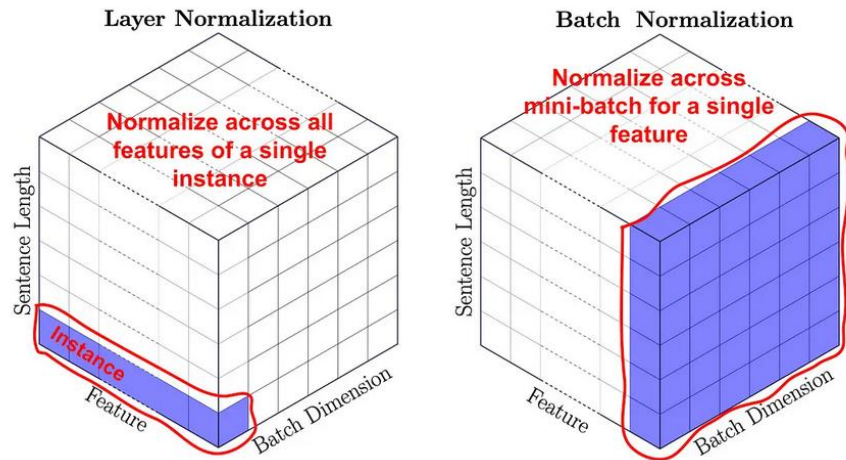
$$\mu^l = \frac{1}{H} \sum_{i=1}^H a_i^l \quad \sigma^l = \sqrt{\frac{1}{H} \sum_{i=1}^H (a_i^l - \mu^l)^2}$$

$$h_i = f\left(\frac{g_i}{\sigma_i} (a_i - \mu_i) + b_i\right)$$

Batch Normalization for RNN

Batch normalization drawbacks:

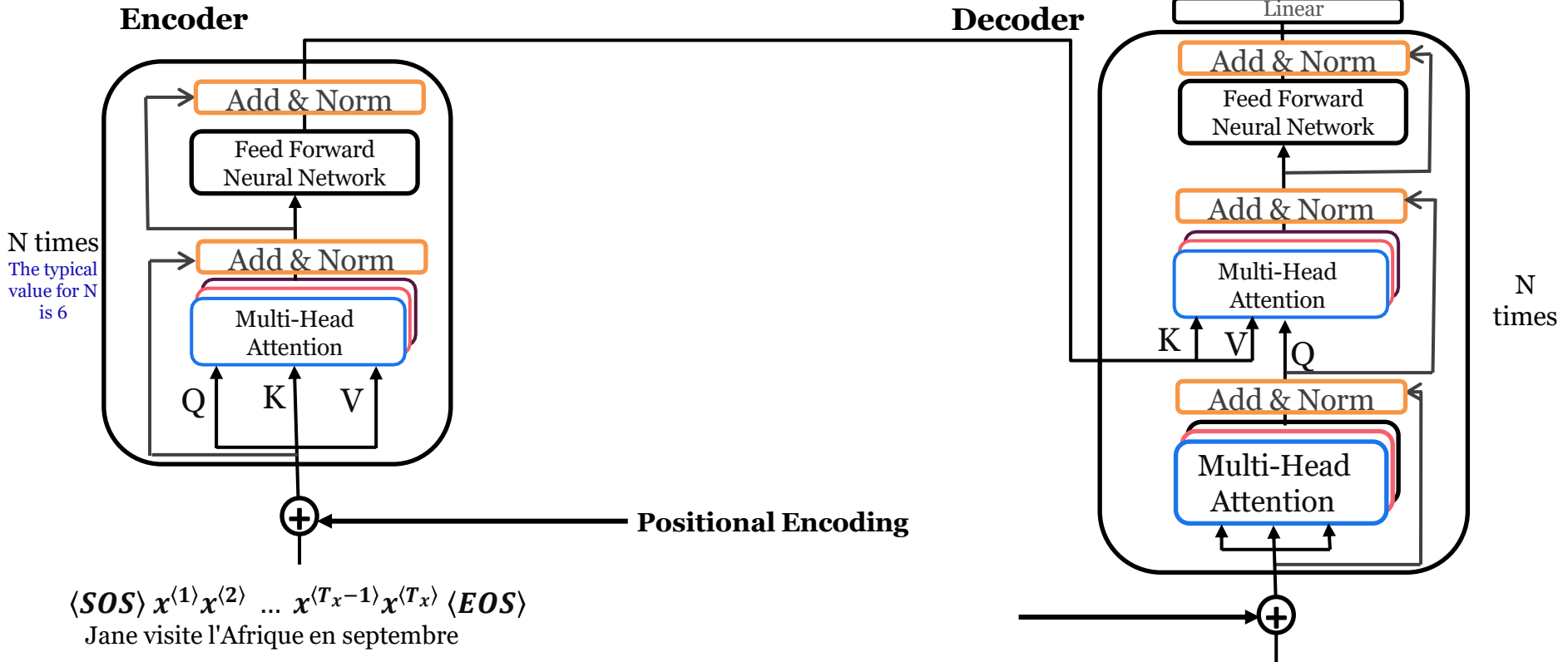
- Applying batch normalization to RNNs appears to require different statistics for different time-steps.
- BN cannot be applied to online learning tasks or to extremely large distributed models where the mini-batches have to be small.



Transformer Details

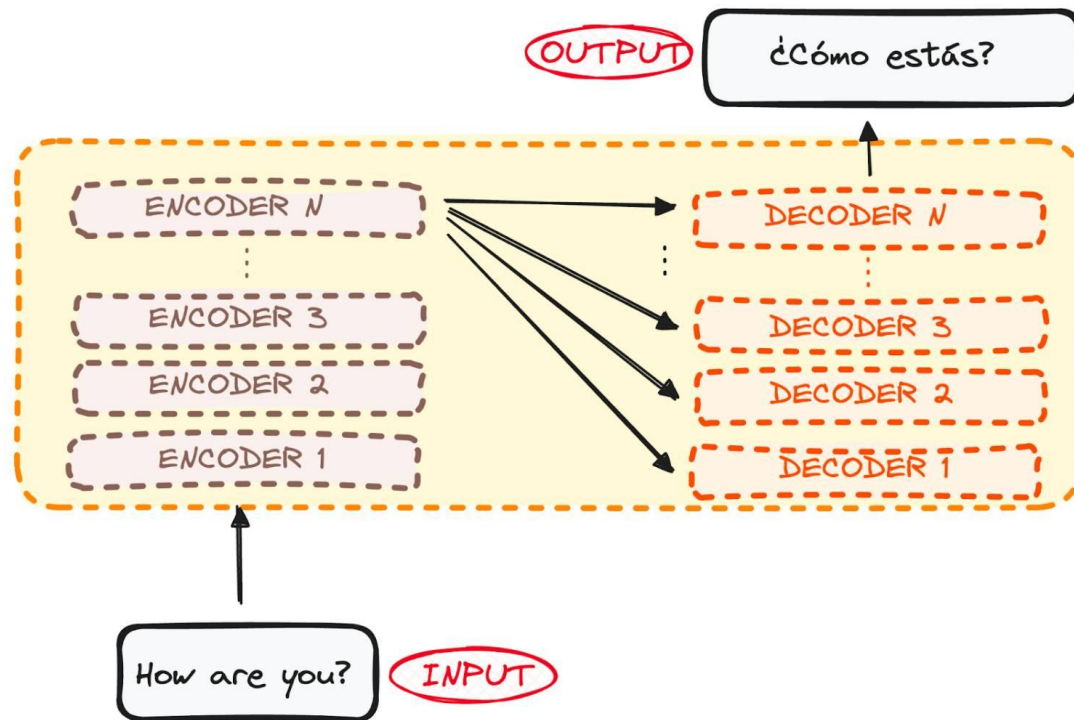
- The output of the encoding block contains contextual semantic embedding and positional encoding information
- Residual connections are used to pass positional information in the transformer architecture.
- “Add & Norm” is a layer similar to batch-norm, helping to speed up learning.
- This “Add & Norm” layer is repeated throughout the architecture.
- Positional encoding and linear layers are used in the transformer.
- A “SoftMax” layer predicts the next word in the decoder block, one word at a time.
- The transformer architecture is designed to handle sequential data.

Transformer Details



$\langle SOS \rangle$ Jane visits Africa in September

Transformer Details



Multi-head attention Dimensions

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$
$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

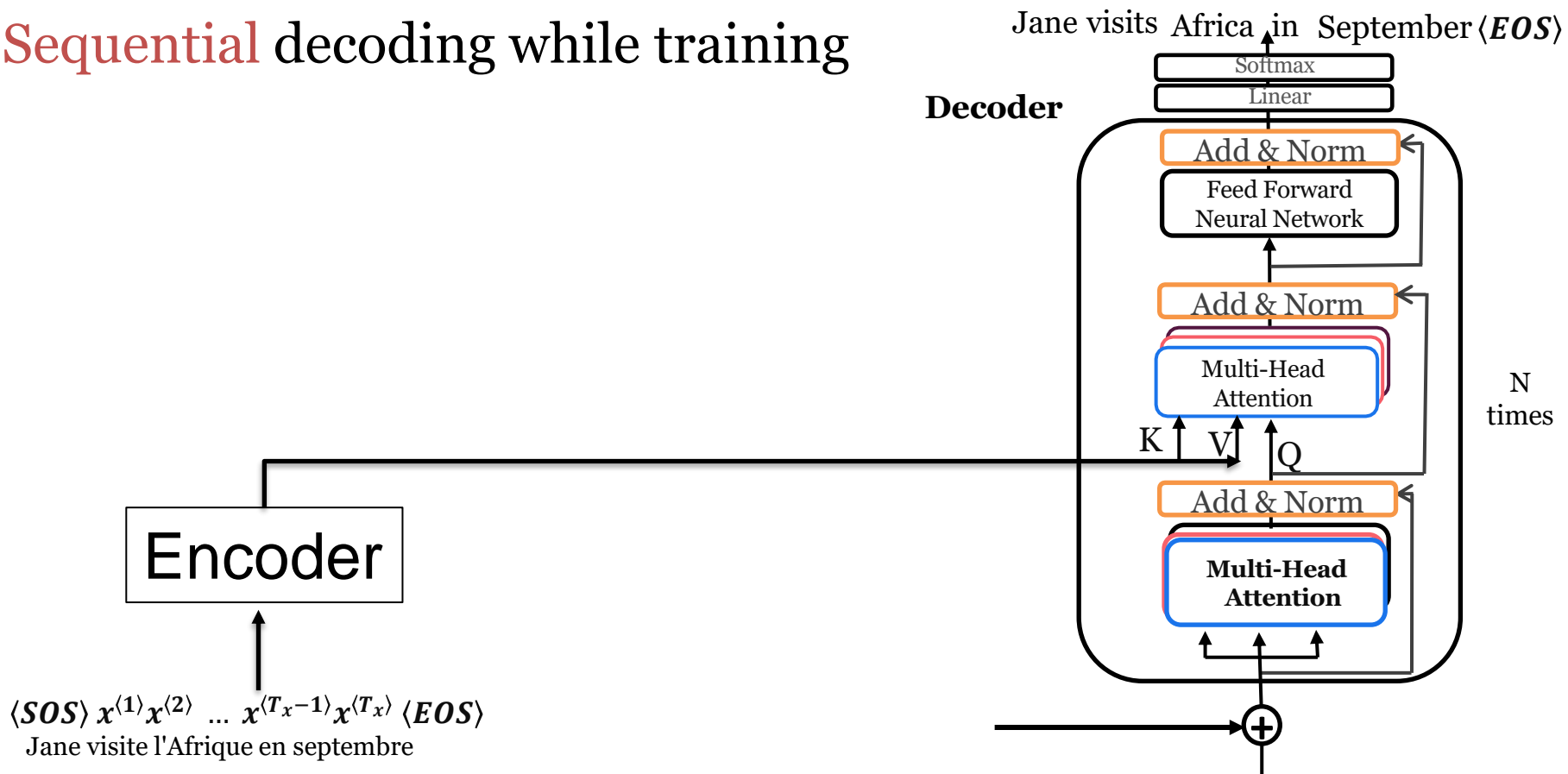
Where the projections are parameter matrices $W_i^Q \in \mathbb{R}^{d_{\text{model}} \times d_k}$, $W_i^K \in \mathbb{R}^{d_{\text{model}} \times d_k}$, $W_i^V \in \mathbb{R}^{d_{\text{model}} \times d_v}$ and $W^O \in \mathbb{R}^{hd_v \times d_{\text{model}}}$.

In this work we employ $h = 8$ parallel attention layers, or heads. For each of these we use $d_k = d_v = d_{\text{model}}/h = 64$. Due to the reduced dimension of each head, the total computational cost is similar to that of single-head attention with full dimensionality.

The dimensionality of input and output is $d_{\text{model}} = 512$,

$$\text{Attention}(QW^Q, KW^K, VW^V) = \text{softmax}\left(\frac{(QW^Q)(KW^K)^T}{\sqrt{d_k}}\right)(VW^V) \quad Q, K, V \in \mathbb{R}^{T \times d_{\text{model}}}$$

Sequential decoding while training

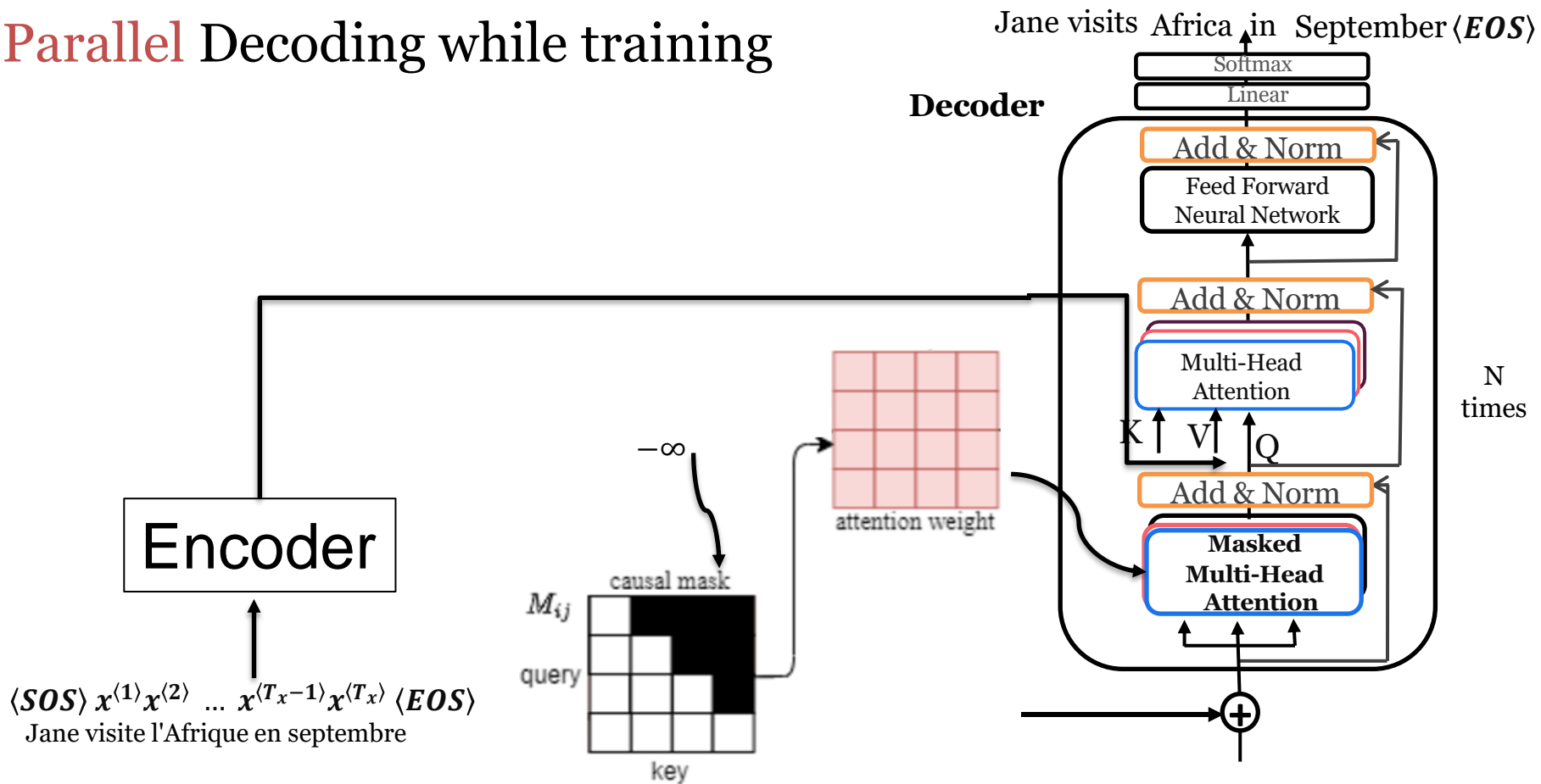


$\langle SOS \rangle$ Jane visits Africa in September $\langle EOS \rangle$

Transformer Details: Masked Multi-Head Attention

- The “**Masked Multi-Head Attention**” is significant during the training process of the transformer network.
- It is used with datasets of correct French to English translations.
- The transformer network predicts one word at a time.
- Understanding how the transformer network is trained is important.
- “Masked Multi-head attention” enhances the transformer's ability to learn translation patterns.
- Masking is used during training to simulate prediction behavior.
- Masking is used during training to mimic test time or prediction scenarios.
- The blocked-out portion of the sentence helps the network predict the next word accurately.
- The aim is to see if the network can perform well given a perfect first part of the translation.

Parallel Decoding while training



$\langle SOS \rangle$ Jane visits Africa in September $\langle EOS \rangle$