



Deep Learning

Dr. Mehran Safayani

safayani@iut.ac.ir

safayani.iut.ac.ir



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



https://github.com/safayani/deep_learning_course

Department of Electrical and computer engineering, Isfahan university of technology, Isfahan, Iran

Attention models & Transformer

Examples of sequence data

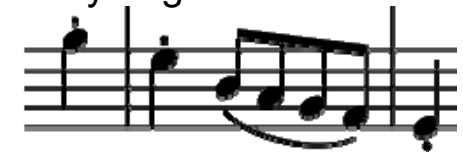
Speech recognition



→ "The quick brown fox jumped over the lazy dog."

Music generation

∅



Sentiment classification

"There is nothing to like in this movie."



DNA sequence analysis

AGCCCCTGTGAGGAACTAG

→ AGCCCCTGTGAGGAACTAG

Machine translation

Voulez-vous chanter avec moi?

→ Do you want to sing with me?

Video activity recognition



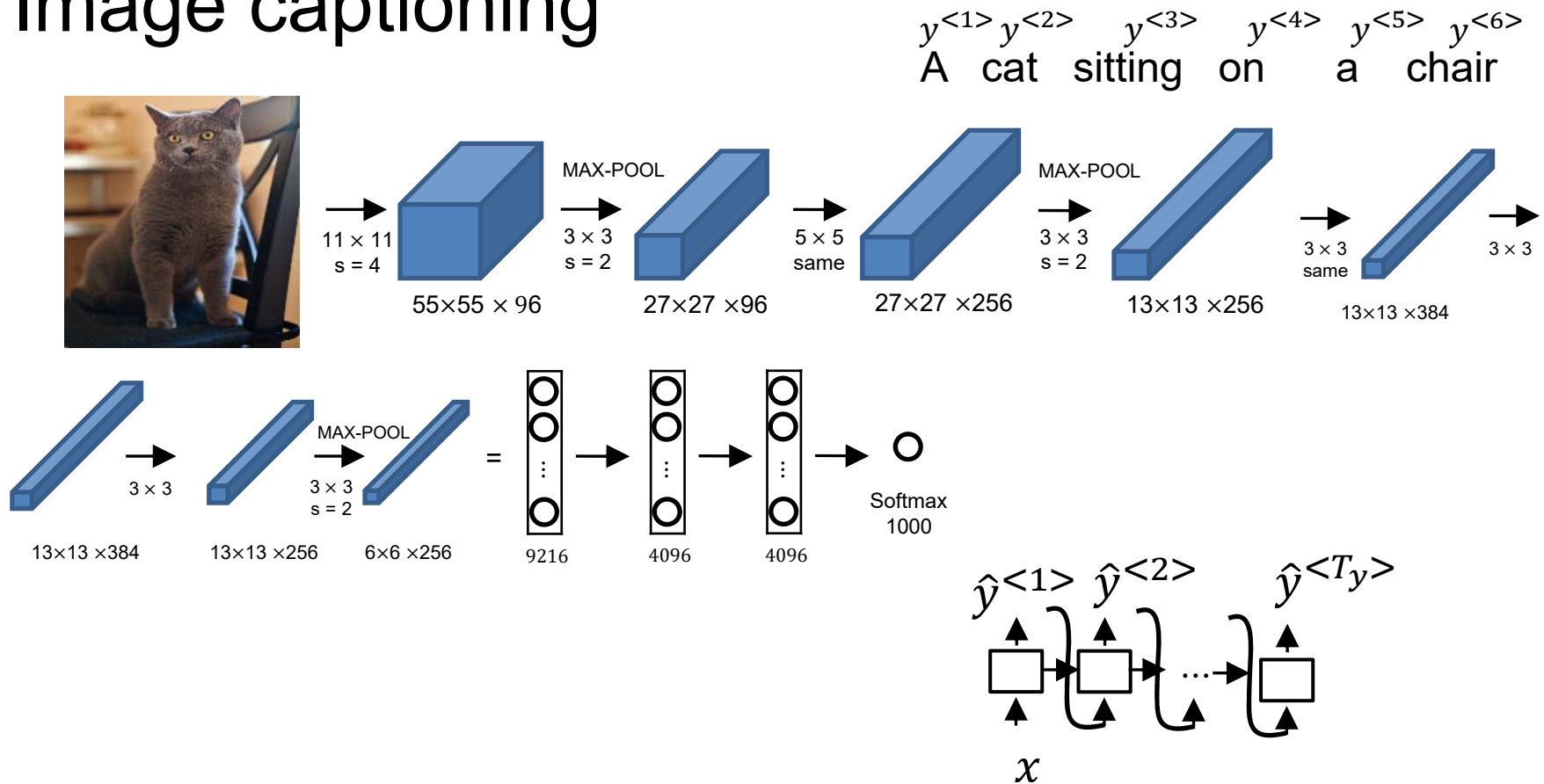
→ Running

Name entity recognition

Yesterday, Harry Potter met Hermione Granger.

→ Yesterday, **Harry Potter** met **Hermione Granger**.

Image captioning



[Mao et. al., 2014. Deep captioning with multimodal recurrent neural networks]

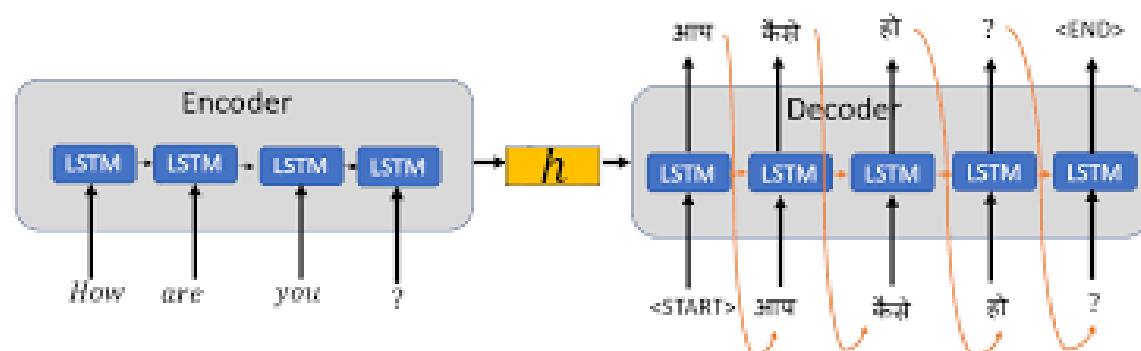
[Vinyals et. al., 2014. Show and tell: Neural image caption generator]

[Karpathy and Li, 2015. Deep visual-semantic alignments for generating image descriptions]

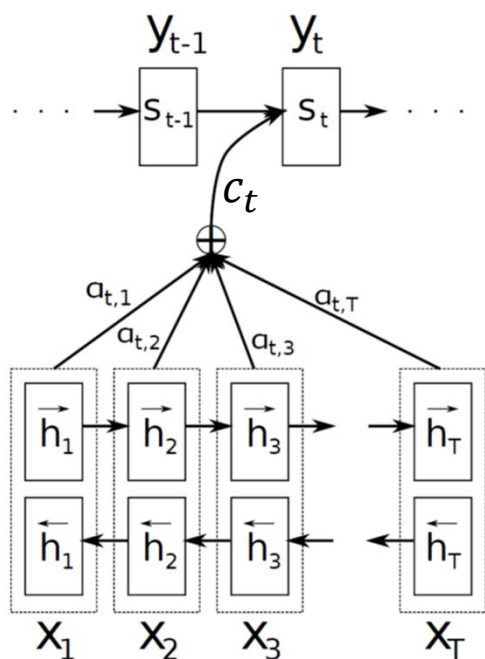
Attention model

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



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

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

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

$$e_{ij} = a(s_{i-1}, h_j) \quad \text{a is a feed forward NN}$$

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

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.

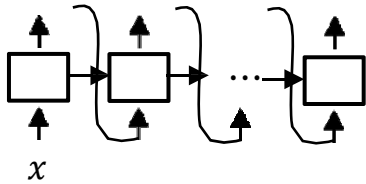
The most important distinguishing feature of this approach from the basic encoder–decoder is that it does not attempt to encode a whole input sentence into a single fixed-length vector. Instead, it encodes the input sentence into a sequence of vectors and chooses a subset of these vectors adaptively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

Transformers Motivation

Increased complexity,
sequential

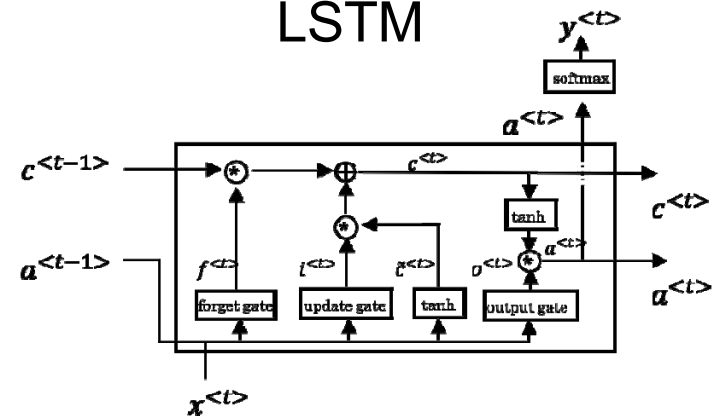


RNN



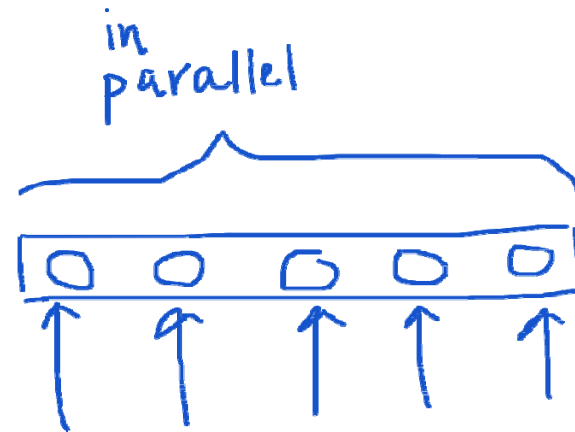
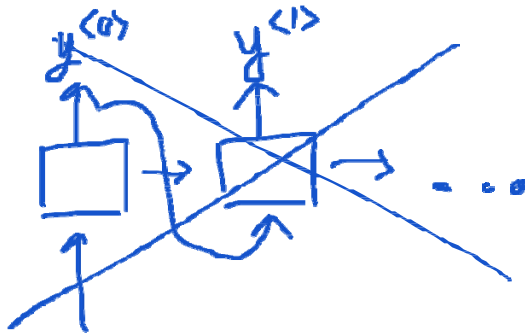
GRU

LSTM



Transformers Intuition

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



[Vaswani et al. 2017, Attention Is All You Need]

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

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

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

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

Transformers Attention

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

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

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

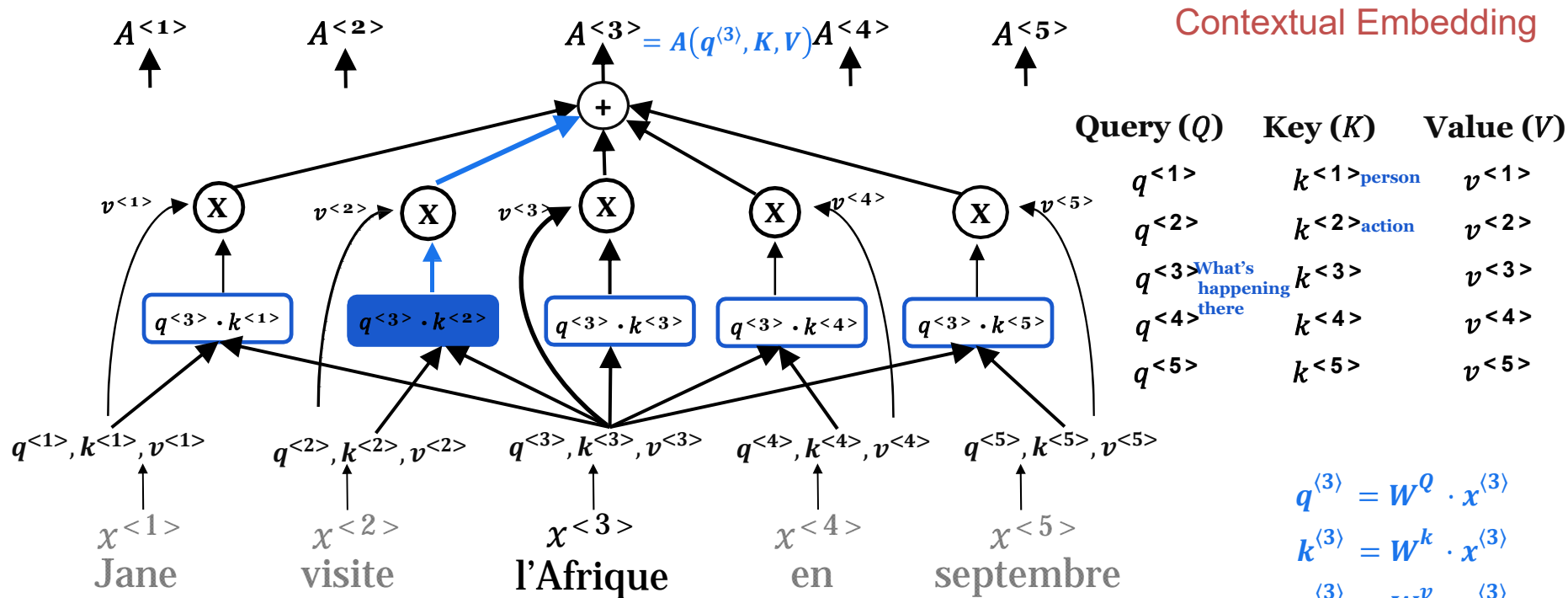
$x^{<5>}$
septembre

Self-Attention

If you put all of these computation together:

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

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



[Vaswani et al. 2017, Attention]

```
# B: Batch size
# T: Sequence length or max token size e.g. 512 for BERT. 'T' because of 'Time steps = Sequence length'
# D: Dimensions of the model embedding vector, which is  $d_{model}$  in the paper.
# H or h: Number of multi attention heads in Multi-head attention
```

```
def calculate_dot_product_similarities(
    query: Tensor,
    key: Tensor,
) -> Tensor:
    """
    Calculate similarity scores between queries and keys using dot product.
```

Args:

```
query: embedding vector of query of shape (B, h, T, d_k)
key: embedding vector of key of shape (B, h, T, d_k)
```

Returns: Similarities (closeness) between q and k of shape (B, h, T, T) where last (T, T) represents relations between all query elements in T sequence against all key elements in T sequence. If T is people in an organization, (T,T) represents all (cartesian product) social connections among them. The relation considers d k number of features.



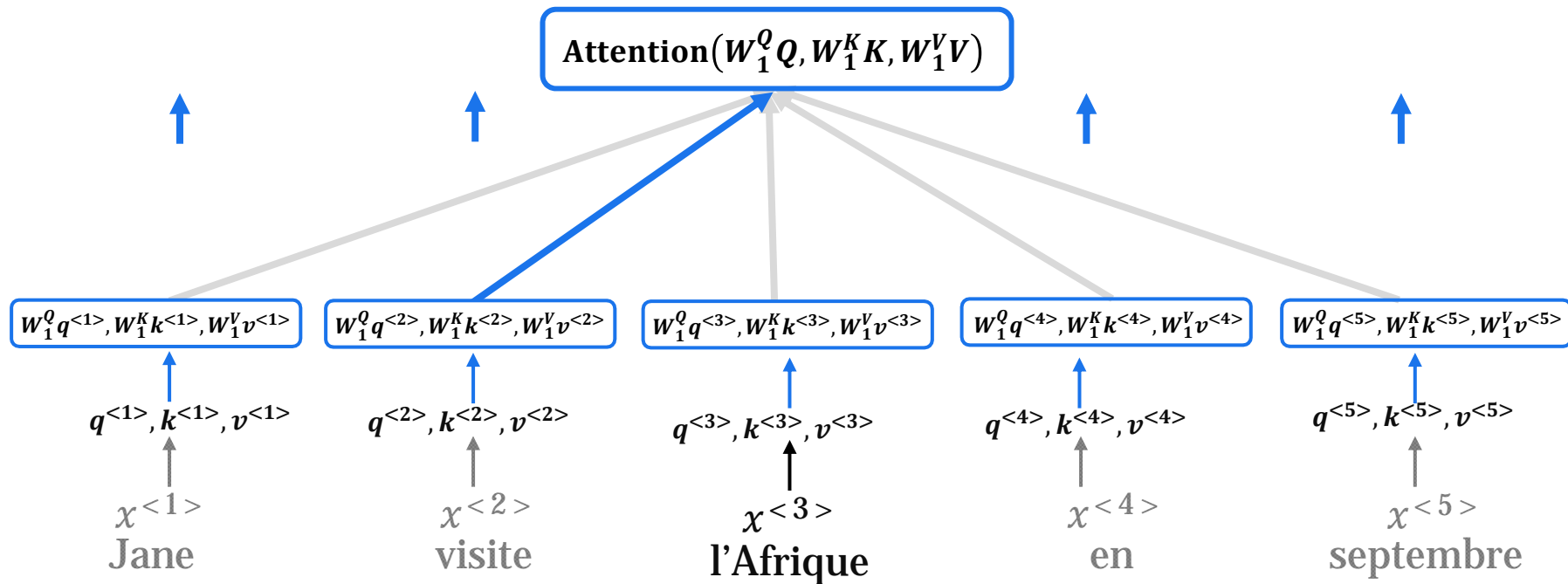
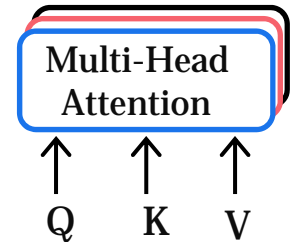
```
# -----  
# Relationship between k and q as the first MatMul using dot product similarity:  
# (B, h, T, d_k) @ (B, hH, d_k, T) ---> (B, h, T, T)  
# -----
```

```
similarities = query @ key.transpose(-2, -1)      # dot product
return similarities                                # shape:(B, h, T, T)
```

Multi-Head Attention

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

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



[Vaswani et al. 2017, Attention Is All You Need]

Multi-Head Attention

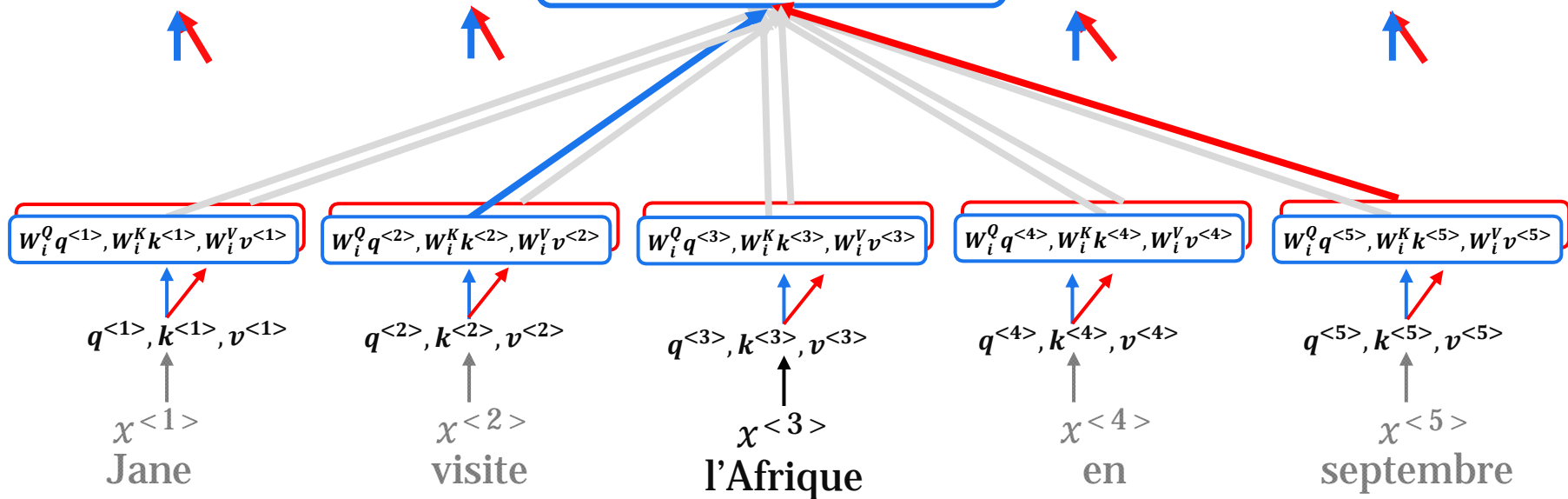
second head: W_2^Q, W_2^K, W_2^V - When is sth happening in Africa?

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

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

Multi-Head
Attention

↑ Q ↑ K ↑ V



[Vaswani et al. 2017, Attention Is All You Need]

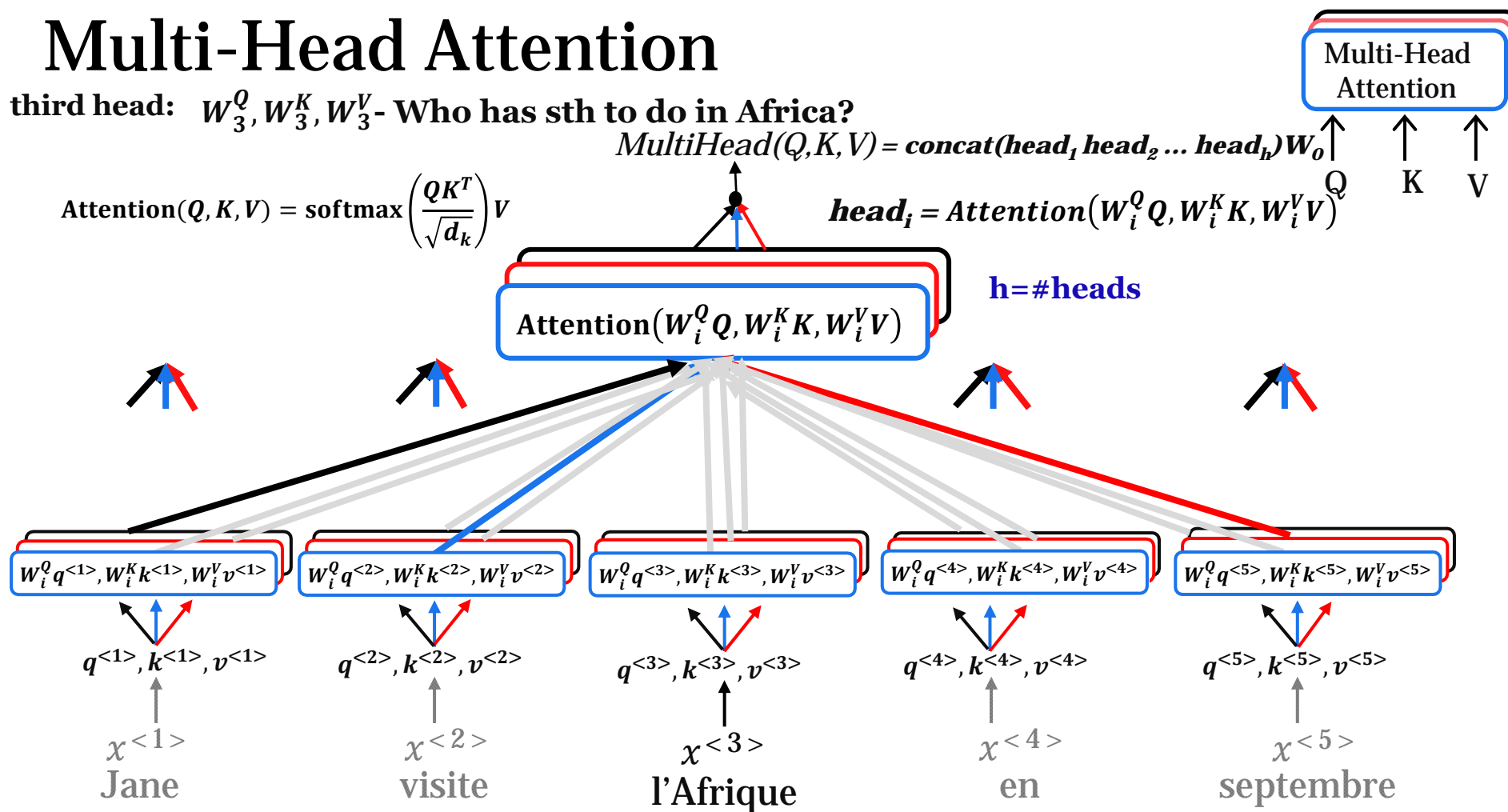
Multi-Head Attention

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

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

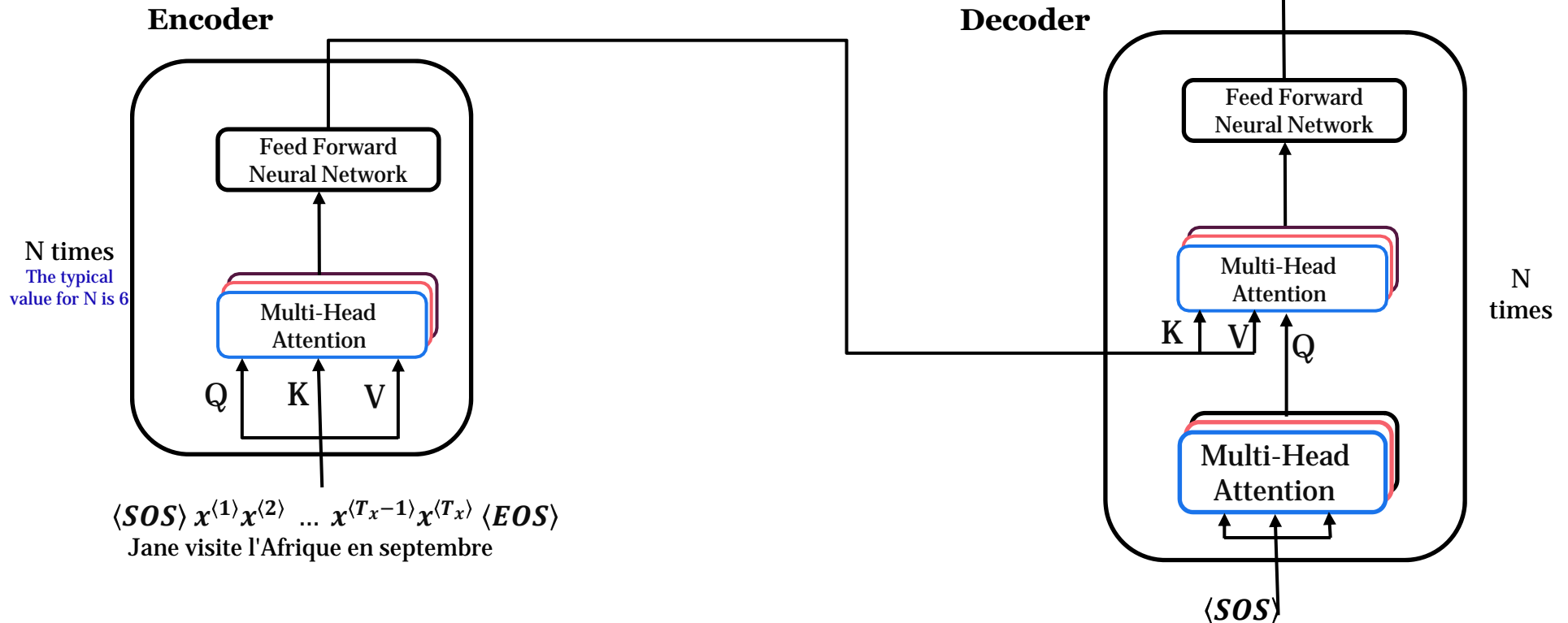
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V$$

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



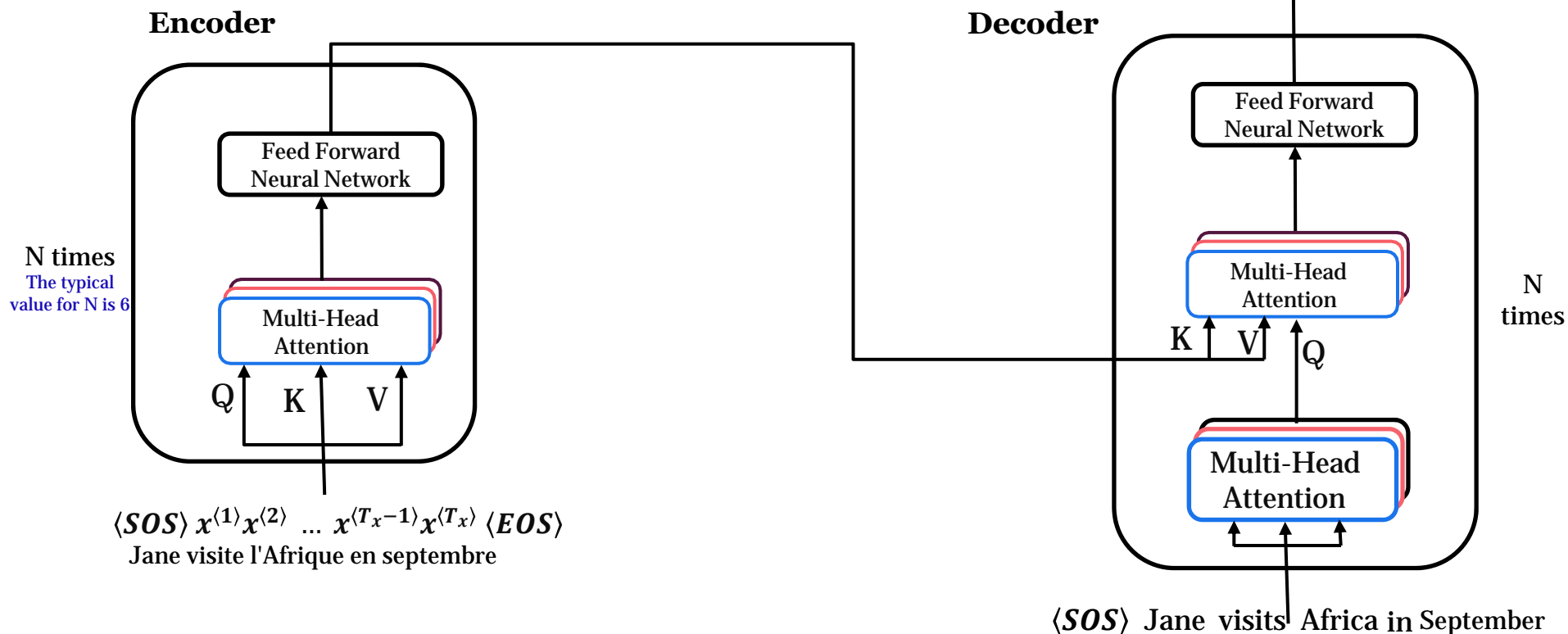
[Vaswani et al. 2017, Attention Is All You Need]

Transformer



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

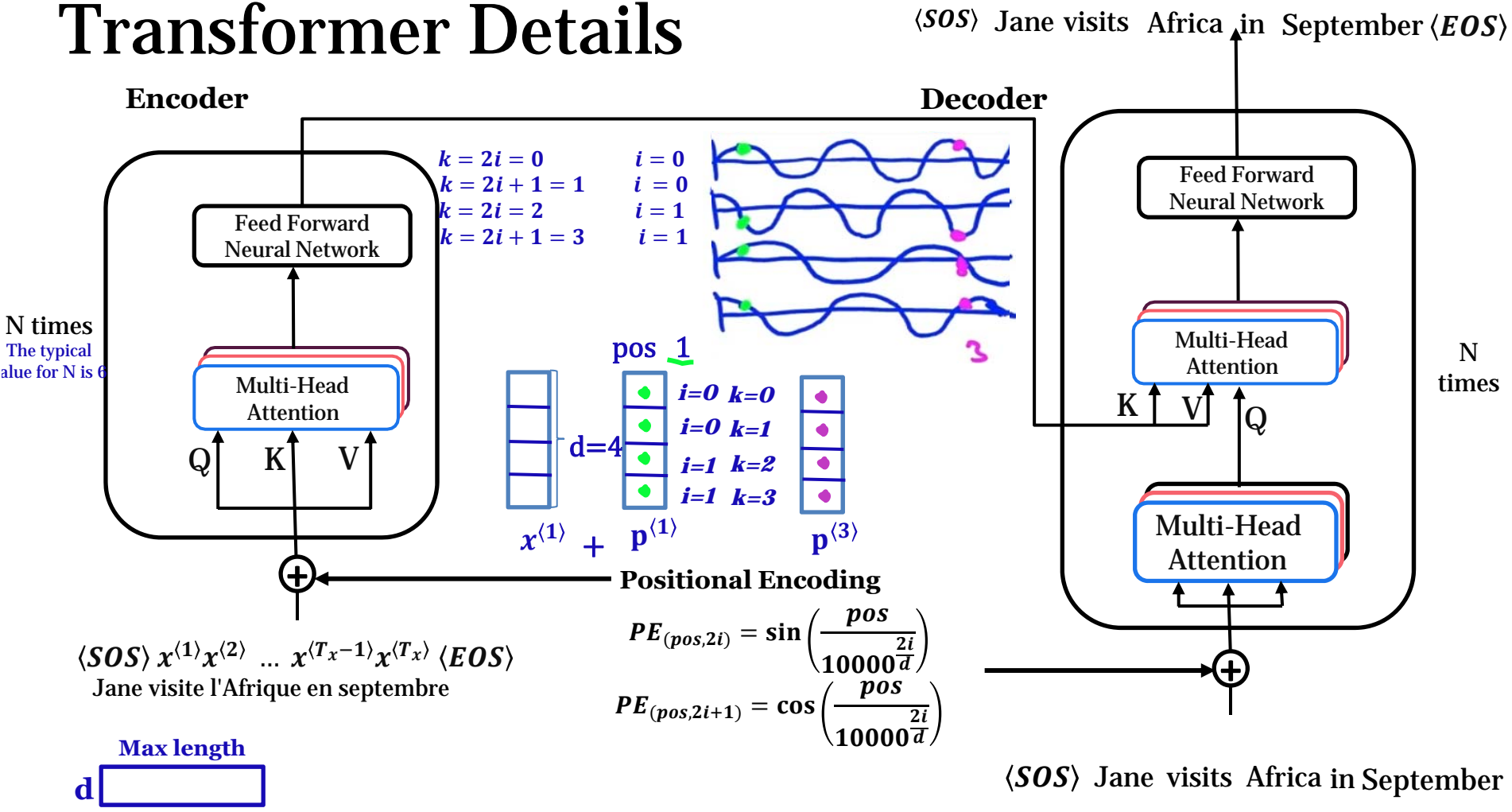
Transformer



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.

Transformer Details



Layer Normalizaion

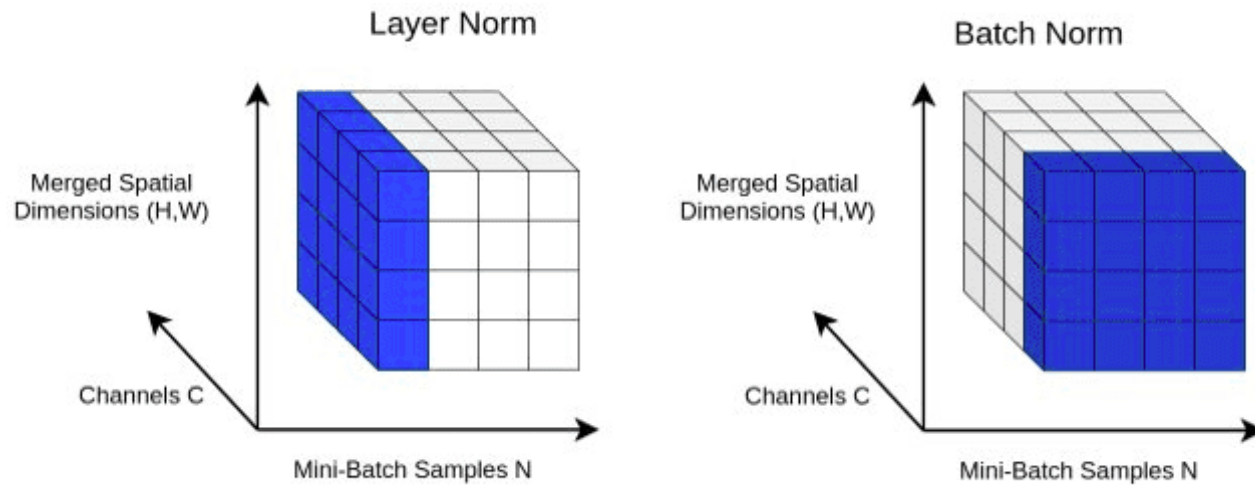
batch normalization drawbacks:

- 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 norm vs Layer norm

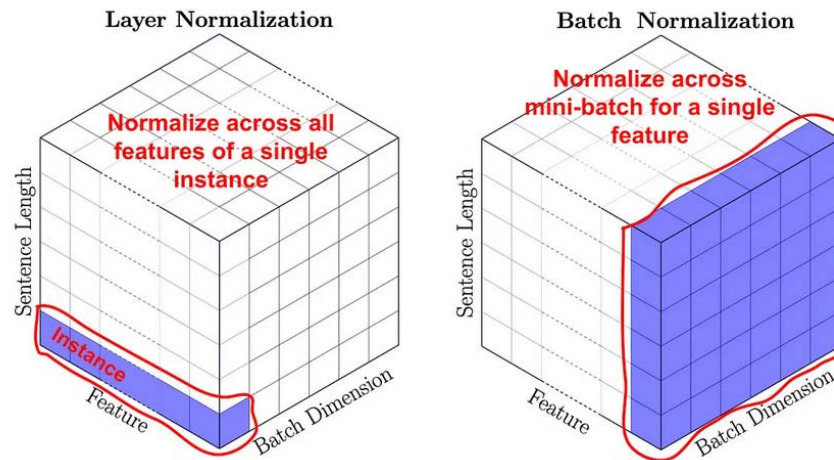


$$y = \frac{x - \mu}{\sqrt{\sigma^2 + \epsilon}} \odot \gamma + \beta$$

Batch Normalization for RNN

Batch normalization drawbacks:

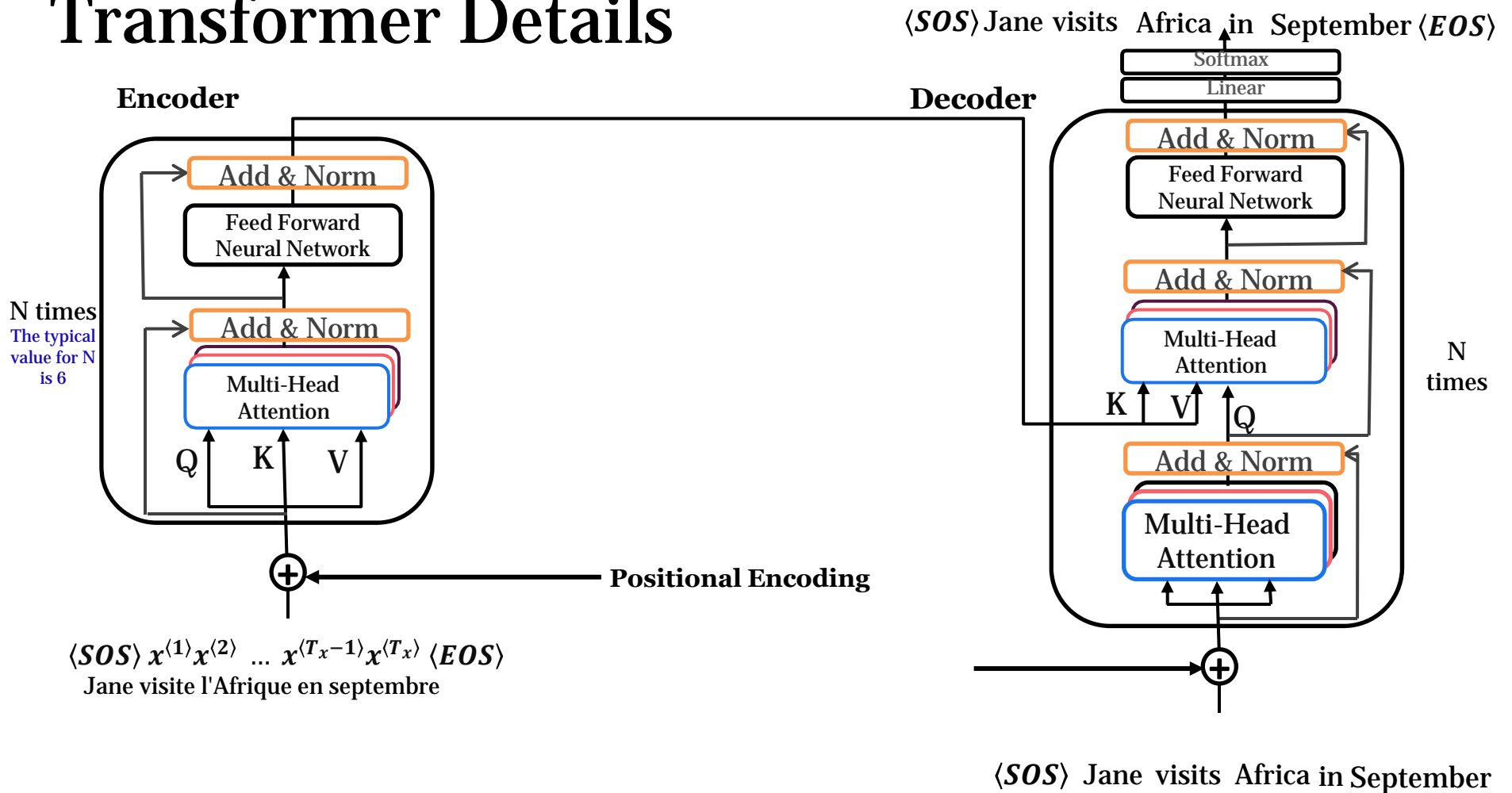
- Applying batch normalization to RNNs appears to require different statistics for different time-steps.
- BN cannot be applied to online learning tasks or to extremely large distributed models where the mini-batches have to be small.



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.

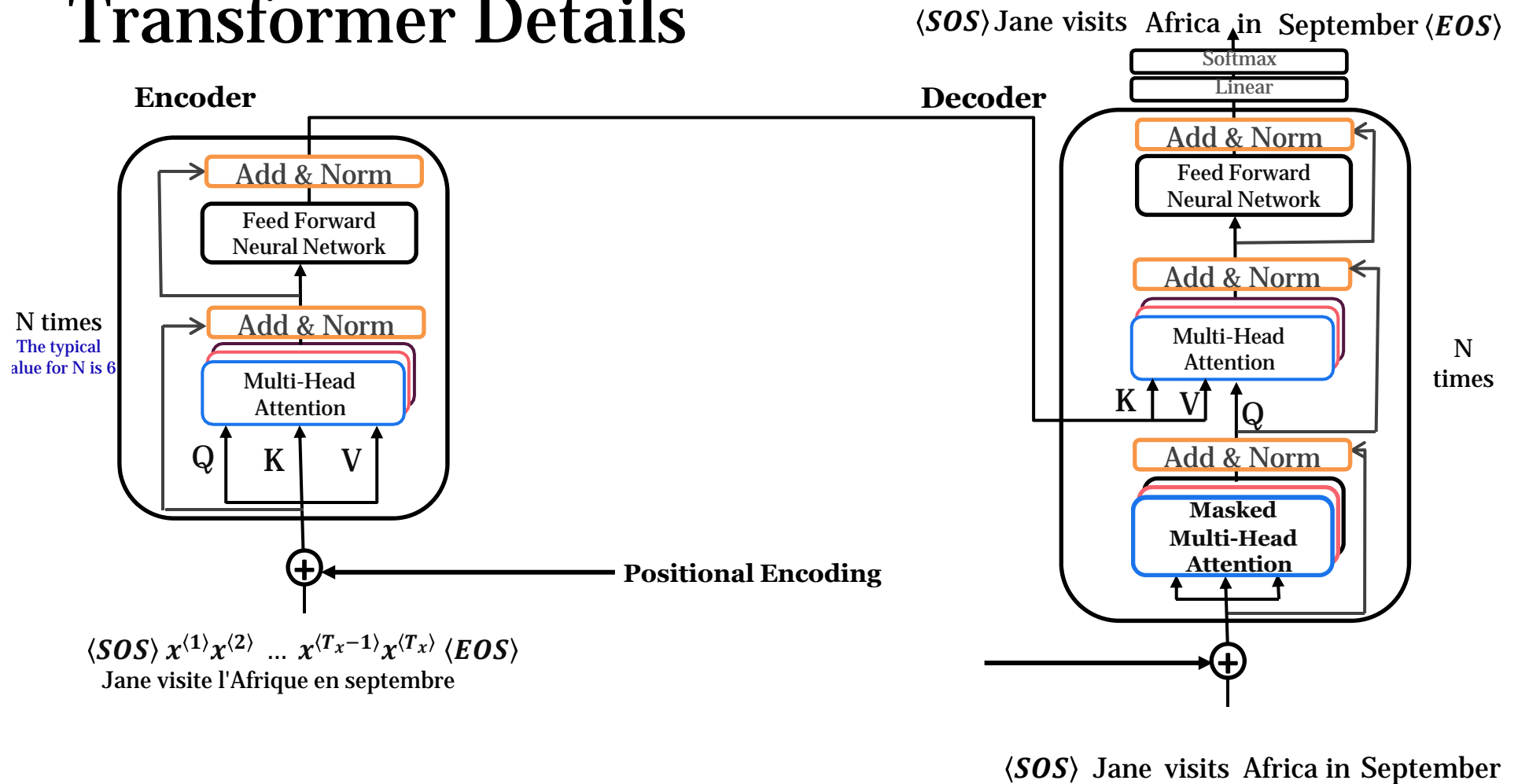
Transformer Details



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.

Transformer Details



reference

[https://medium.com/@hunter-j-phillips/layer-normalization-e9ae93eb3c9c#:~:text=In%20natural%20language%20processing%2C%20layer,will%20be%20\(5%2C\).](https://medium.com/@hunter-j-phillips/layer-normalization-e9ae93eb3c9c#:~:text=In%20natural%20language%20processing%2C%20layer,will%20be%20(5%2C).)

<https://medium.com/@bhavtoshrath.umn/batch-normalization-vs-layer-normalization-c76bb3cbf388>