

Copyright Notice

These slides are distributed under the Creative Commons License.

[DeepLearning.AI](#) makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite [DeepLearning.AI](#) as the source of the slides.

For the rest of the details of the license, see <https://creativecommons.org/licenses/by-sa/2.0/legalcode>



deeplearning.ai

Sequence to sequence models

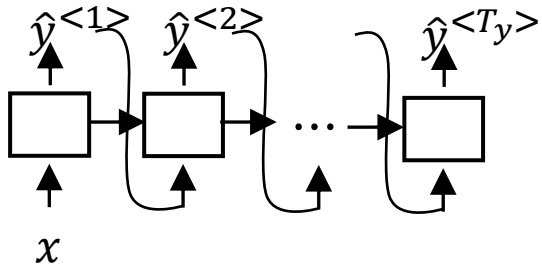
Transformers Intuition

Transformers Motivation

Increased complexity,
sequential

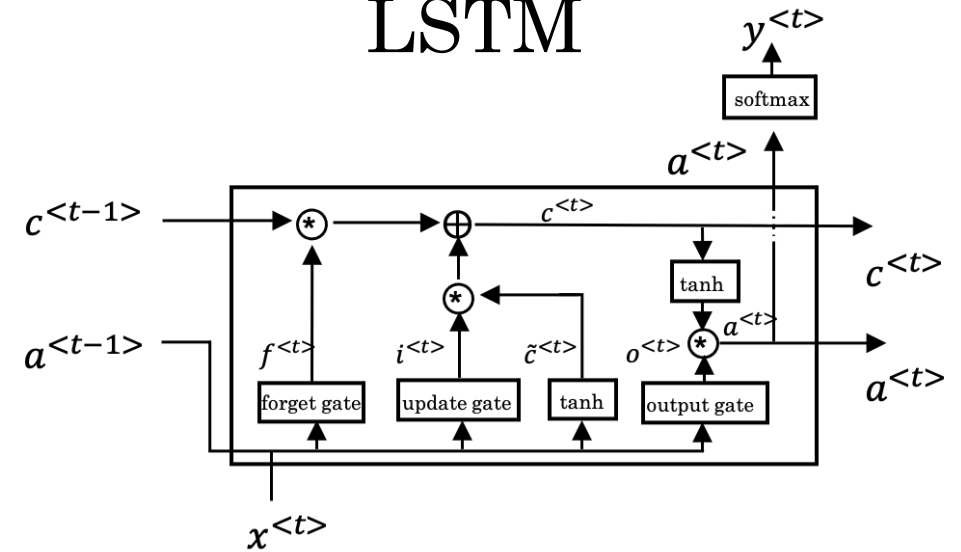


RNN



GRU

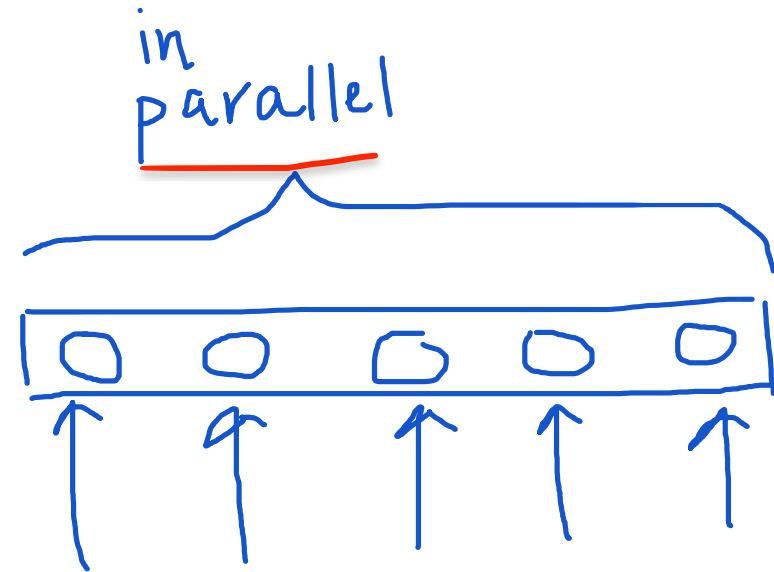
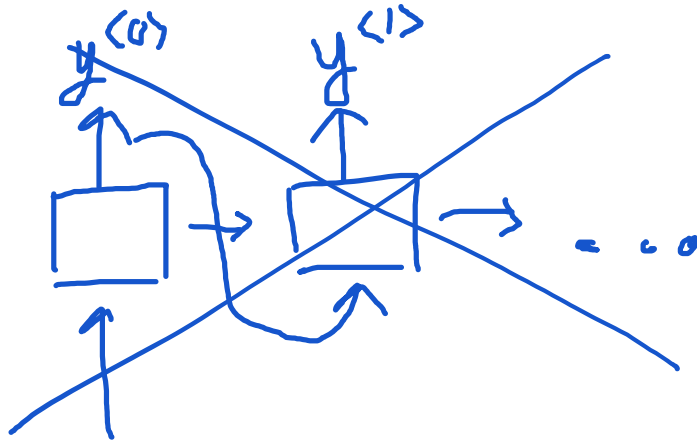
LSTM



Transformers Intuition

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

for loop of self attention





deeplearning.ai

Sequence to sequence models

Self-Attention

Self-Attention Intuition

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

→ calculate for each word

RNN Attention

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

Transformers Attention

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

“word embedding” that consider neighbours words

$x^{<1>}$

Jane

$x^{<2>}$

visite

$x^{<3>}$

l’Afrique

$x^{<4>}$

en

$x^{<5>}$

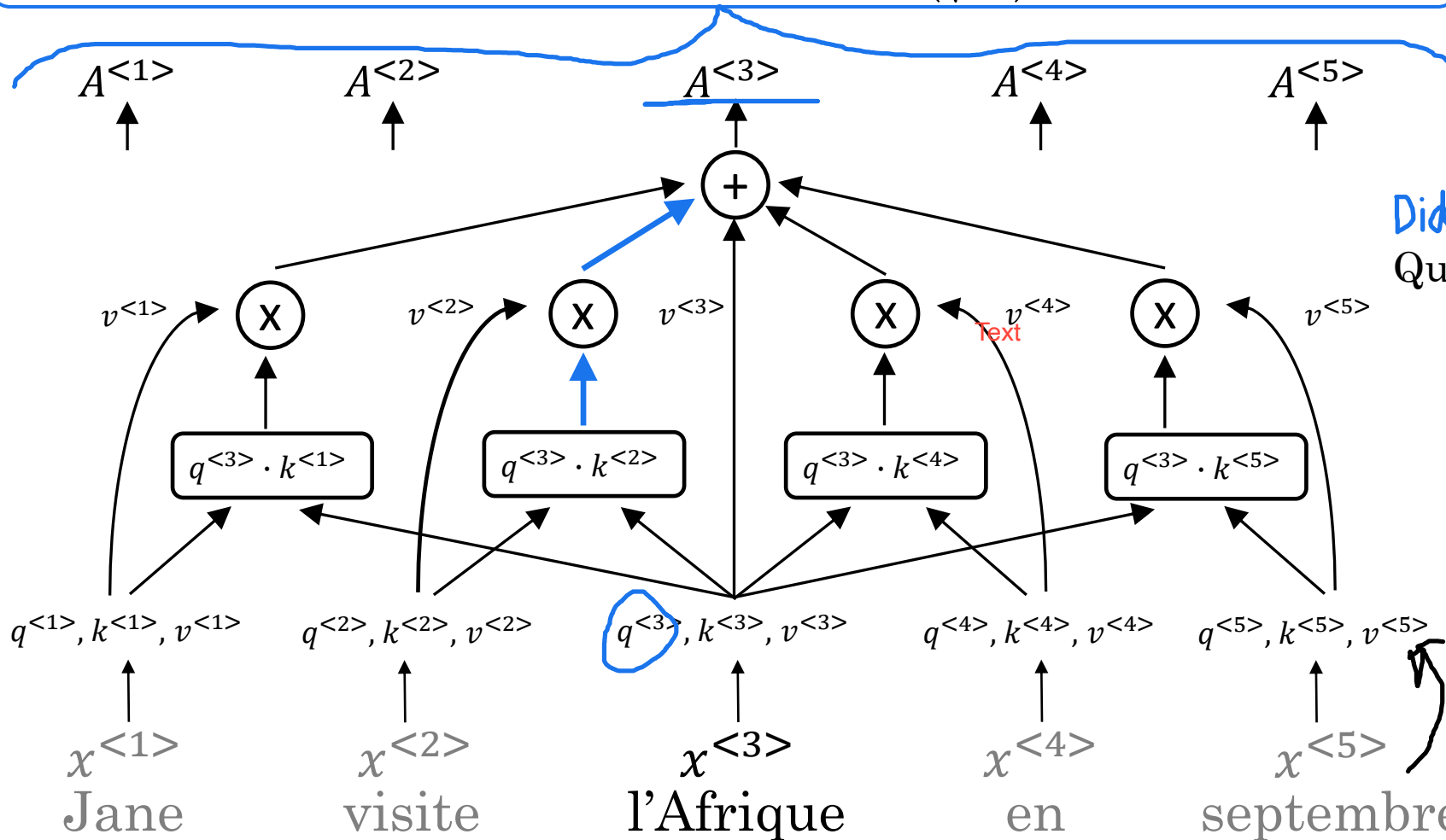
septembre

Self-Attention

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

softmax

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



Did what?

Query (Q)

Key (K)

Value (V)

$q^{<1>}$

$k^{<1>}$

$v^{<1>}$

$q^{<2>}$

$k^{<2>}$

$v^{<2>}$

$q^{<3>}$

$k^{<3>}$

$v^{<3>}$

$q^{<4>}$

$k^{<4>}$

$v^{<4>}$

$q^{<5>}$

$k^{<5>}$

$v^{<5>}$

What's happening there?

person
action
visit

W^Q, W^K, W^V



deeplearning.ai

Sequence to sequence models

Multi-Head Attention

Note that in the previous video you learned that you calculate q , k and v (circled red) by multiplying x with matrices W . In case of the multi-head attention, you don't need to do this, as you already have the matrices W

in each head, and you would effectively do the calculation twice if you did the multiplication here also.

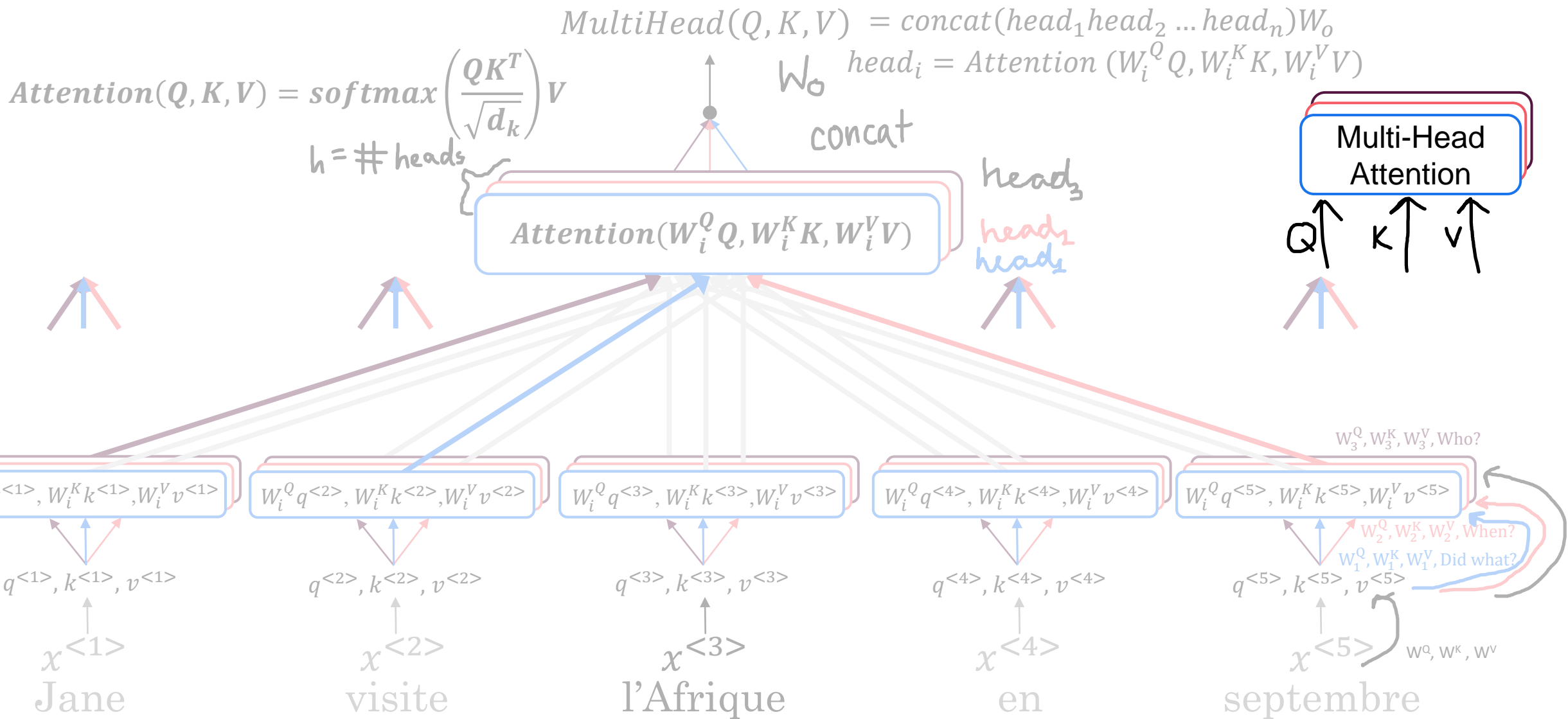
In the simplest case of multi-headed self-attention you would actually use

$q=k=v=x$. The reason we anyway show q , k and v here as different values is that in one part of the transformer

(where you calculate the attention between the input and output) they are not all the same,

as they carry different information as you will learn in the next video.

Multi-Head Attention





deeplearning.ai

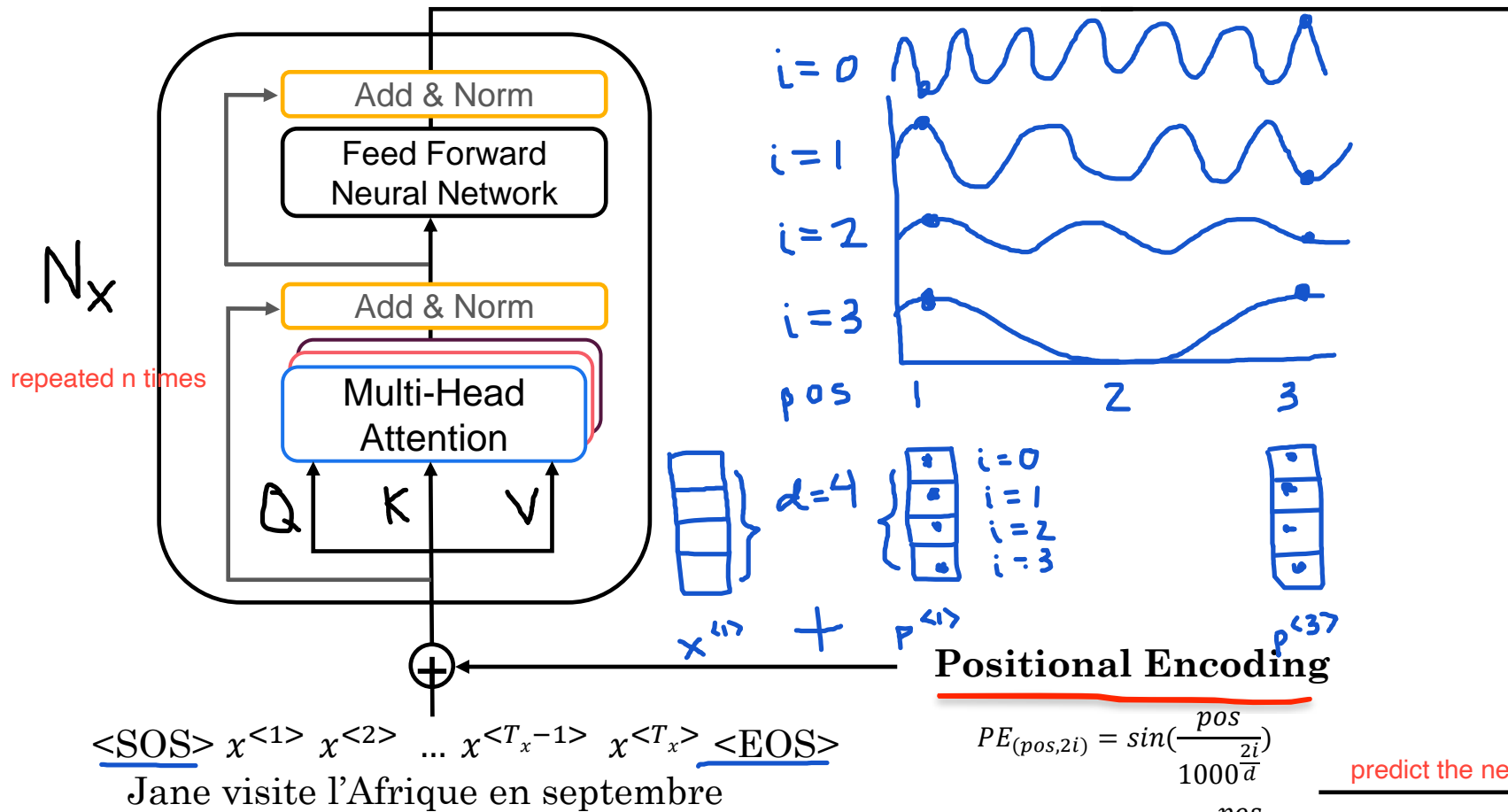
Sequence to sequence models

Transformers

Transformer Details

<SOS> Jane visits Africa in September <EOS>

Encoder



Decoder

