

# 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

---

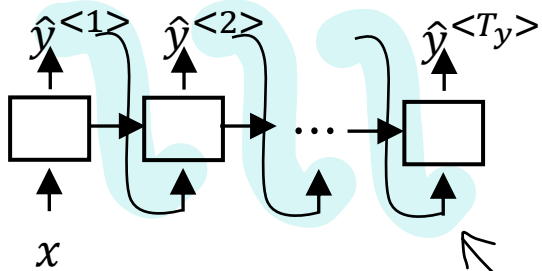
Transformers  
Intuition

# Transformers Motivation

Increased complexity,  
sequential

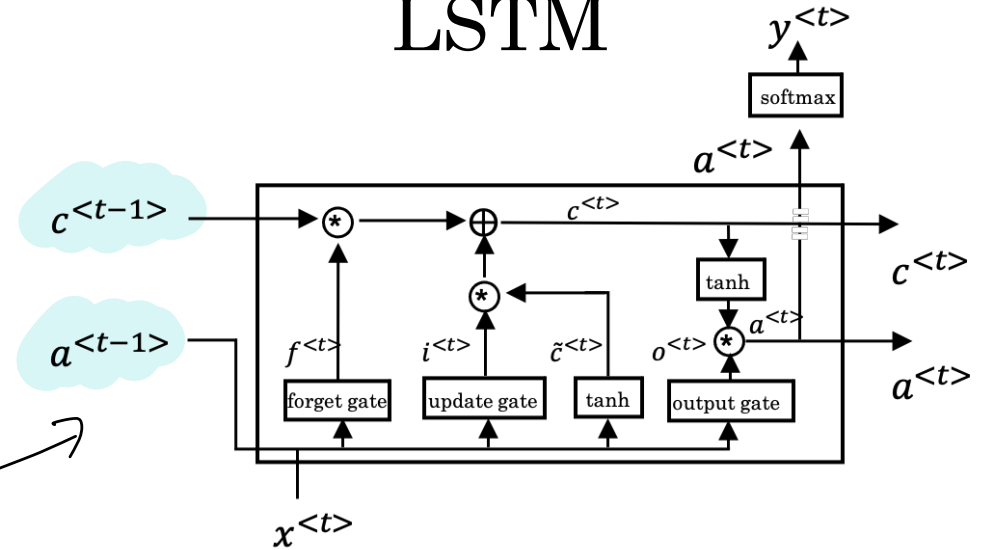


RNN



GRU

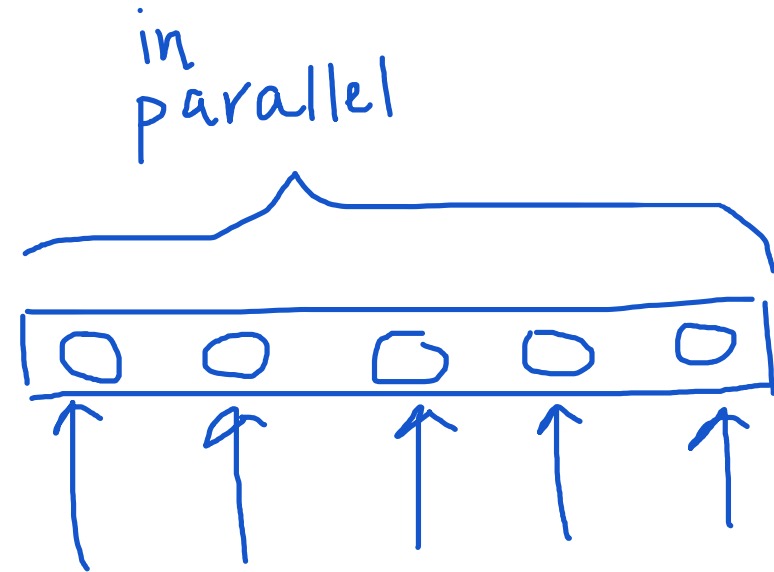
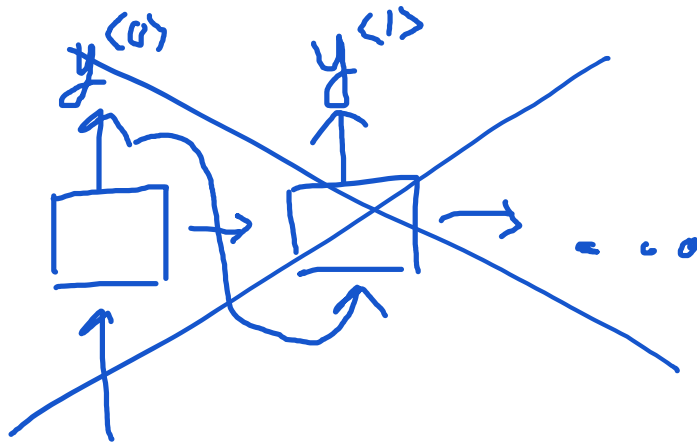
LSTM



are sequential — need to wait for result from previous (&next) timestep

# Transformers Intuition

- Attention + CNN
  - Self-Attention
  - Multi-Head Attention for loop over self-attention





deeplearning.ai

# Sequence to sequence models

---

## Self-Attention

# Self-Attention Intuition

↑ is found using context (surrounding words)

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

↪ calculate for each word in parallel

## RNN Attention

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

## Transformers Attention

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

$x^{<1>}$   
Jane

$x^{<2>}$   
visite

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

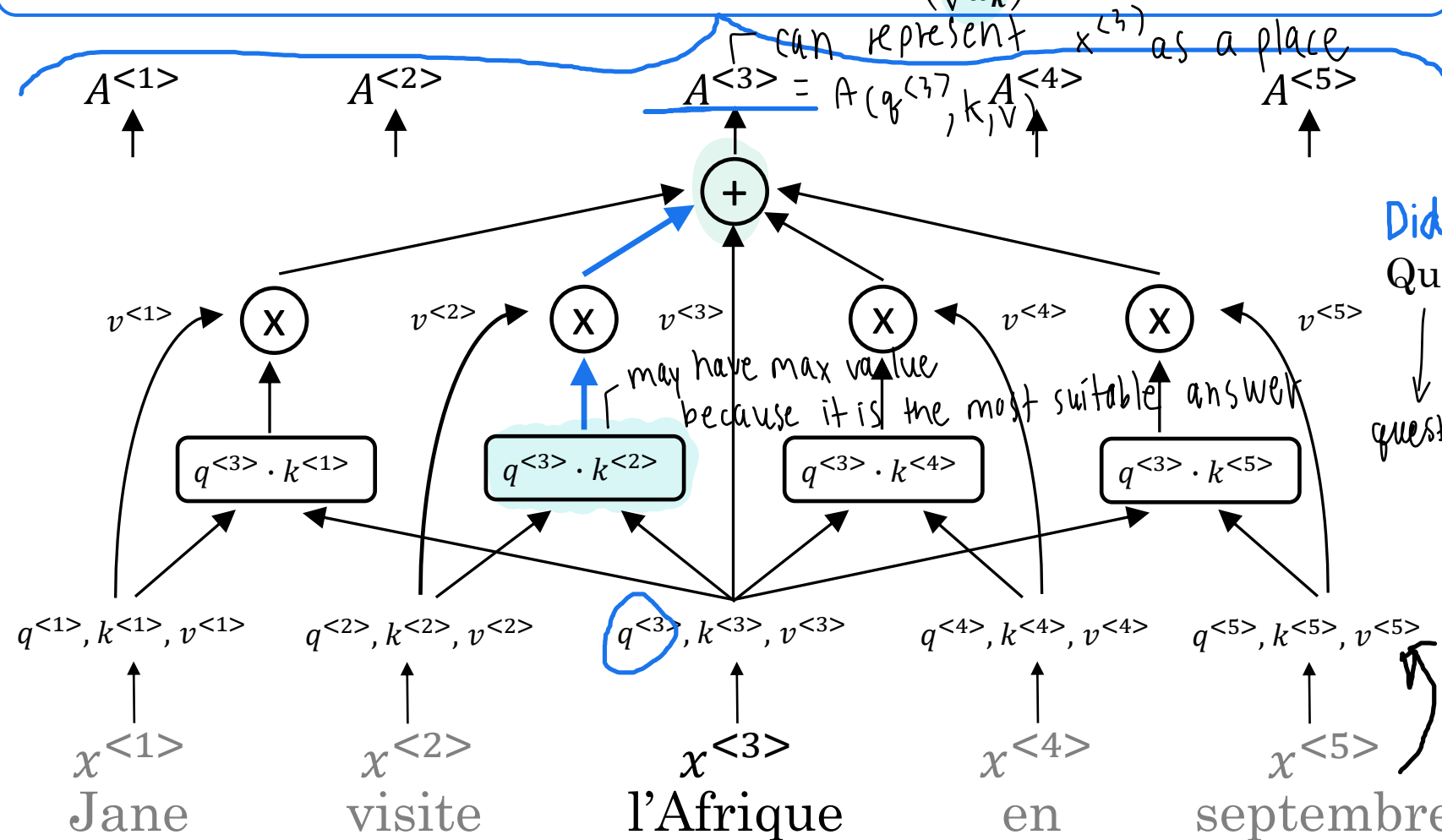
$x^{<4>}$   
en

$x^{<5>}$   
septembre

# Self-Attention

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

for scaling



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

softmax

are learnable

$$q^{(i)} = W^Q \cdot x^{(i)}$$

$$k^{(i)} = W^K \cdot x^{(i)}$$

$$v^{(i)} = W^V \cdot x^{(i)}$$

Did what?

Query (Q)

$q^{(1)}$   
 $q^{(2)}$   
 $q^{(3)}$  question  
 $q^{(4)}$  What's happening there?  
 $q^{(5)}$

Key (K)

$k^{(1)}$  person  
 $k^{(2)}$  action  
 $k^{(3)}$   
 $k^{(4)}$   
 $k^{(5)}$

Value (V)

$v^{(1)}$  Jane  
 $v^{(2)}$  visit  
 $v^{(3)}$   
 $v^{(4)}$   
 $v^{(5)}$

$W^Q, W^K, W^V$



deeplearning.ai

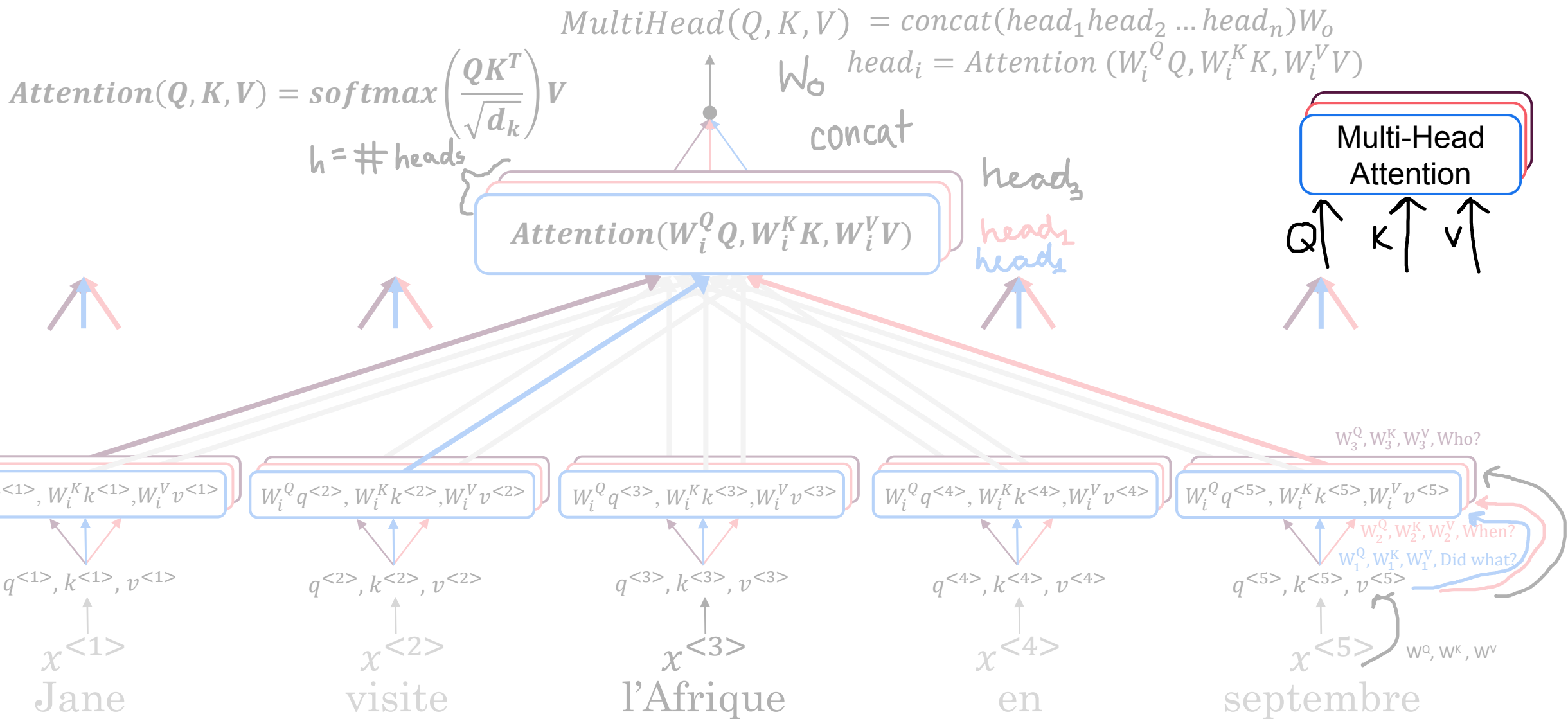
# Sequence to sequence models

---

## Multi-Head Attention



# Multi-Head Attention





deeplearning.ai

# Sequence to sequence models

---

# Transformers

# Transformer Details

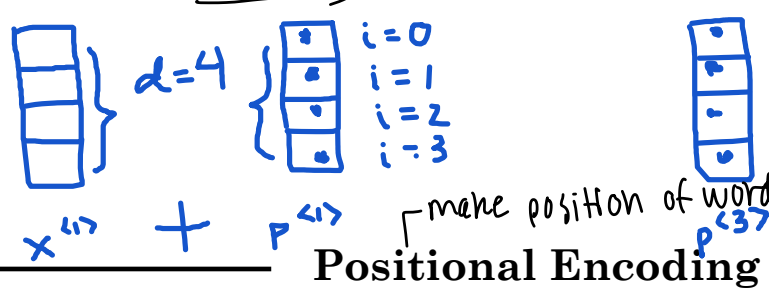
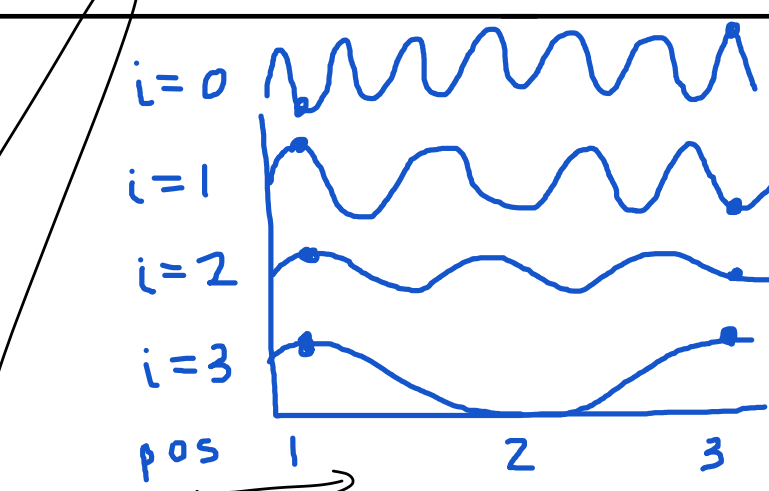
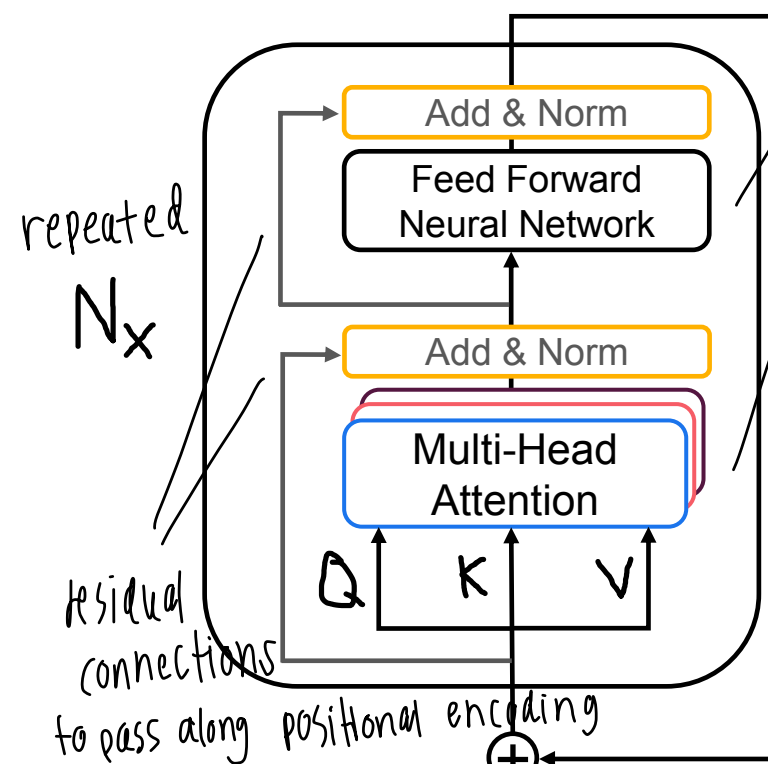
masking = hiding part of input

is predicted from  $\langle \text{SOS} \rangle$  (previous word) &  $K, V$  from encoder  
 $\langle \text{SOS} \rangle$  Jane visits Africa in September  $\langle \text{EOS} \rangle$

## Encoder

output shape =  $d \times \text{max-sequence-length}$

## Decoder

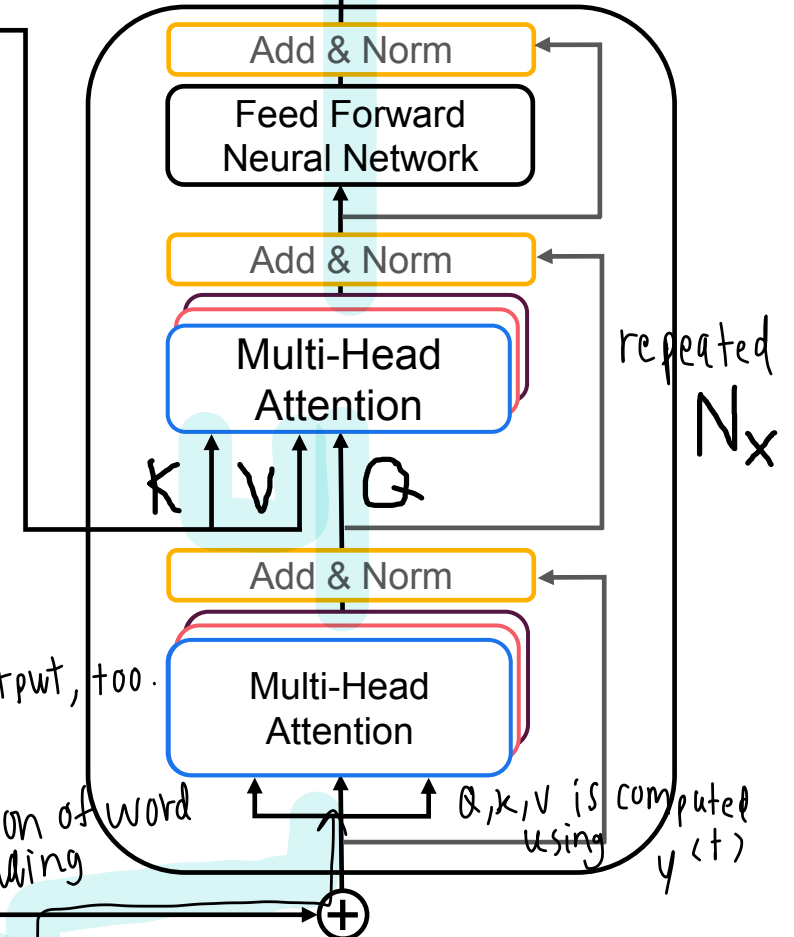


$$PE_{(pos, 2i)} = \sin\left(\frac{pos}{1000^{\frac{2i}{d}}}\right)$$

$$PE_{(pos, 2i+1)} = \cos\left(\frac{pos}{1000^{\frac{2i}{d}}}\right)$$

$i$  starts at zero  
 $pos$  = position of word

$\langle \text{SOS} \rangle x^{(1)} x^{(2)} \dots x^{(T_x-1)} x^{(T_x)} \langle \text{EOS} \rangle$   
 Jane visite l'Afrique en septembre



$\langle \text{SOS} \rangle y^{(1)} y^{(2)} \dots y^{(T_y-1)} y^{(T_y)}$   
 $\langle \text{SOS} \rangle$  Jane visits Africa in September

DONE