# Attention, Transformer and BERT

Prof. Kuan-Ting Lai

2020/6/16
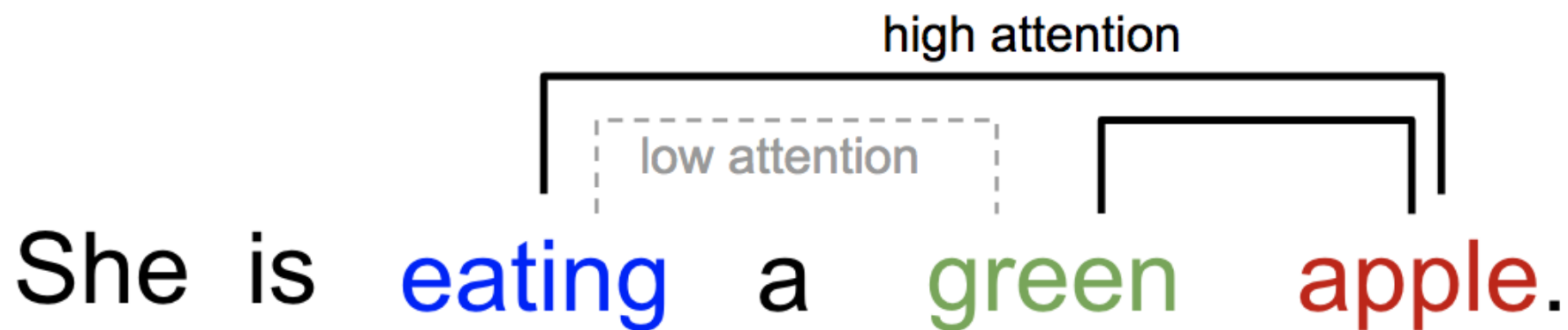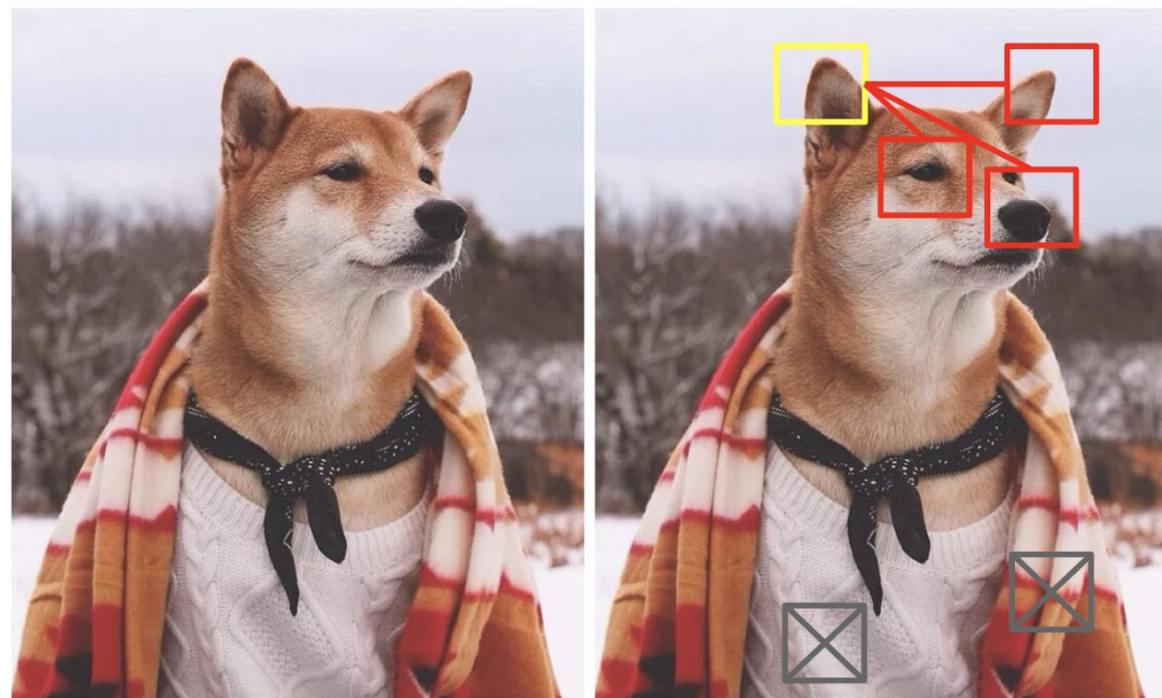
# Attention is All You Need!

A. Waswani et al., *NIPS*, 2017

Google Brain & University of Toronto
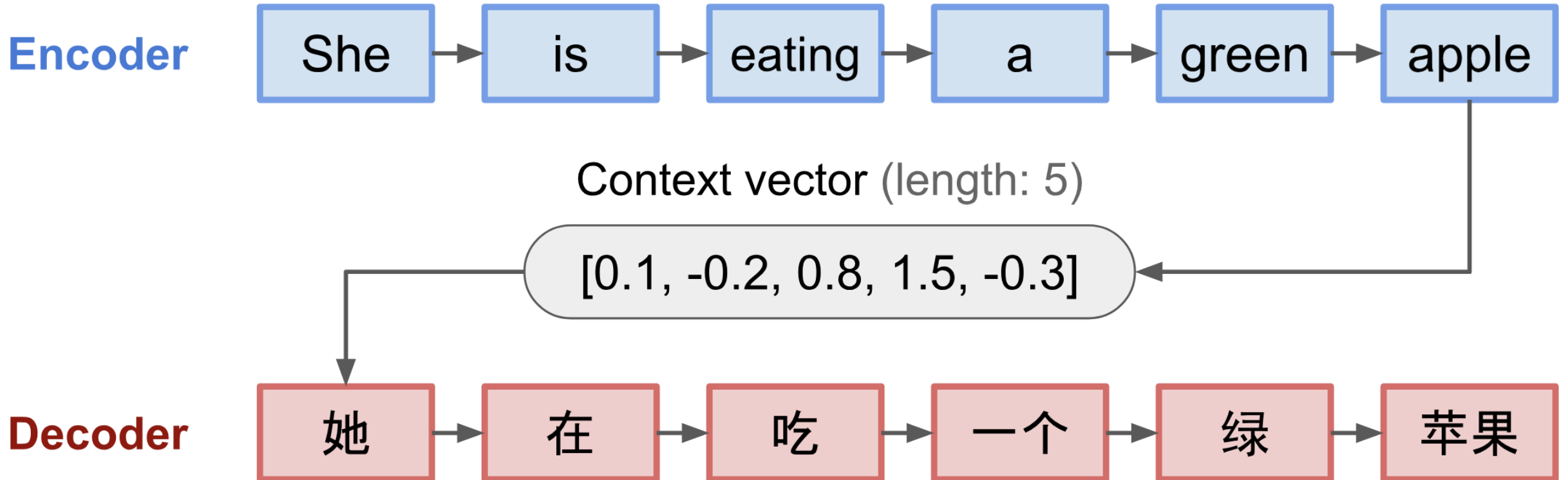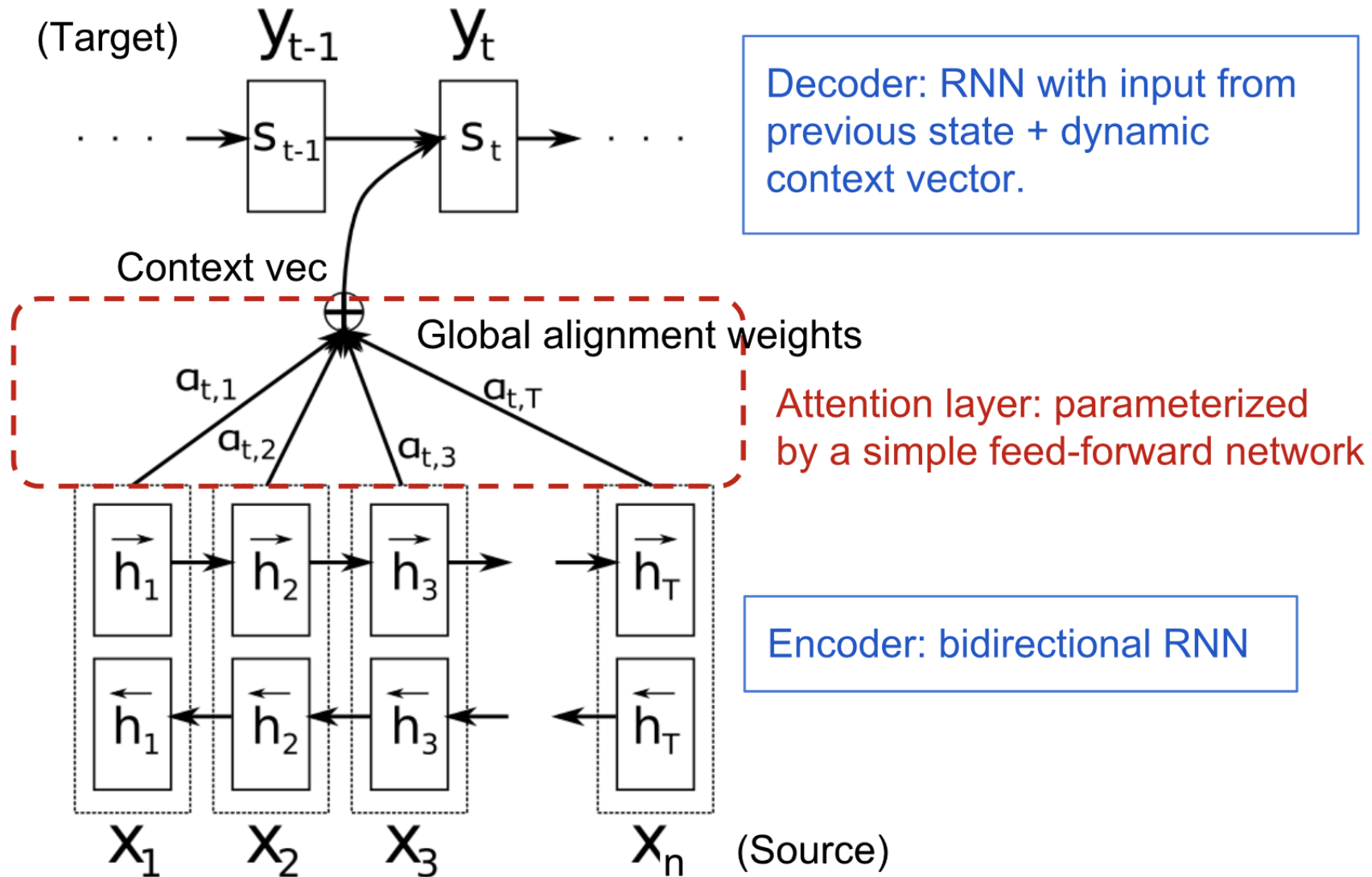
# Attention

- Visual attention and textual attention

# Seq2seq model

- Language translation



**Encoder** | She → is → eating → a → green → apple

Context vector (length: 5)

[0.1, -0.2, 0.8, 1.5, -0.3]

**Decoder** | 她 → 在 → 吃 → 一个 → 绿 → 苹果

# Attention = Vector of Importance Weights



(Target)

$y_{t-1}$   $y_t$

$S_{t-1}$   $S_t$

Context vec

Global alignment weights

$a_{t,1}$   $a_{t,2}$   $a_{t,3}$   $a_{t,T}$

$\overrightarrow{h_1}$   $\overrightarrow{h_2}$   $\overrightarrow{h_3}$   $\overrightarrow{h_T}$

$\overleftarrow{h_1}$   $\overleftarrow{h_2}$   $\overleftarrow{h_3}$   $\overleftarrow{h_T}$

$X_1$   $X_2$   $X_3$   $X_n$   (Source)

Decoder: RNN with input from previous state + dynamic context vector.

Attention layer: parameterized by a simple feed-forward network

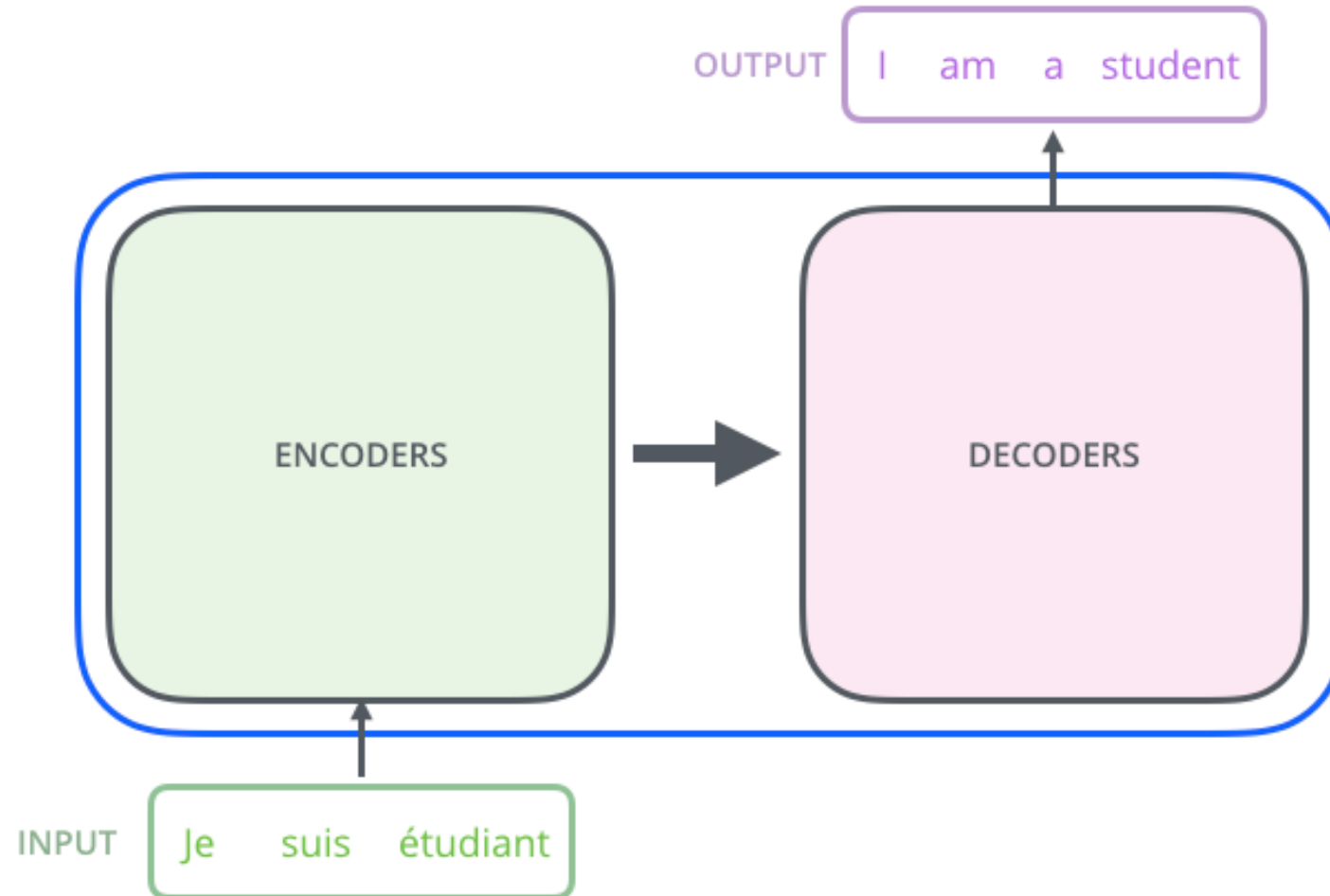**Additive Attention**

Encoder: bidirectional RNN

# Transformer

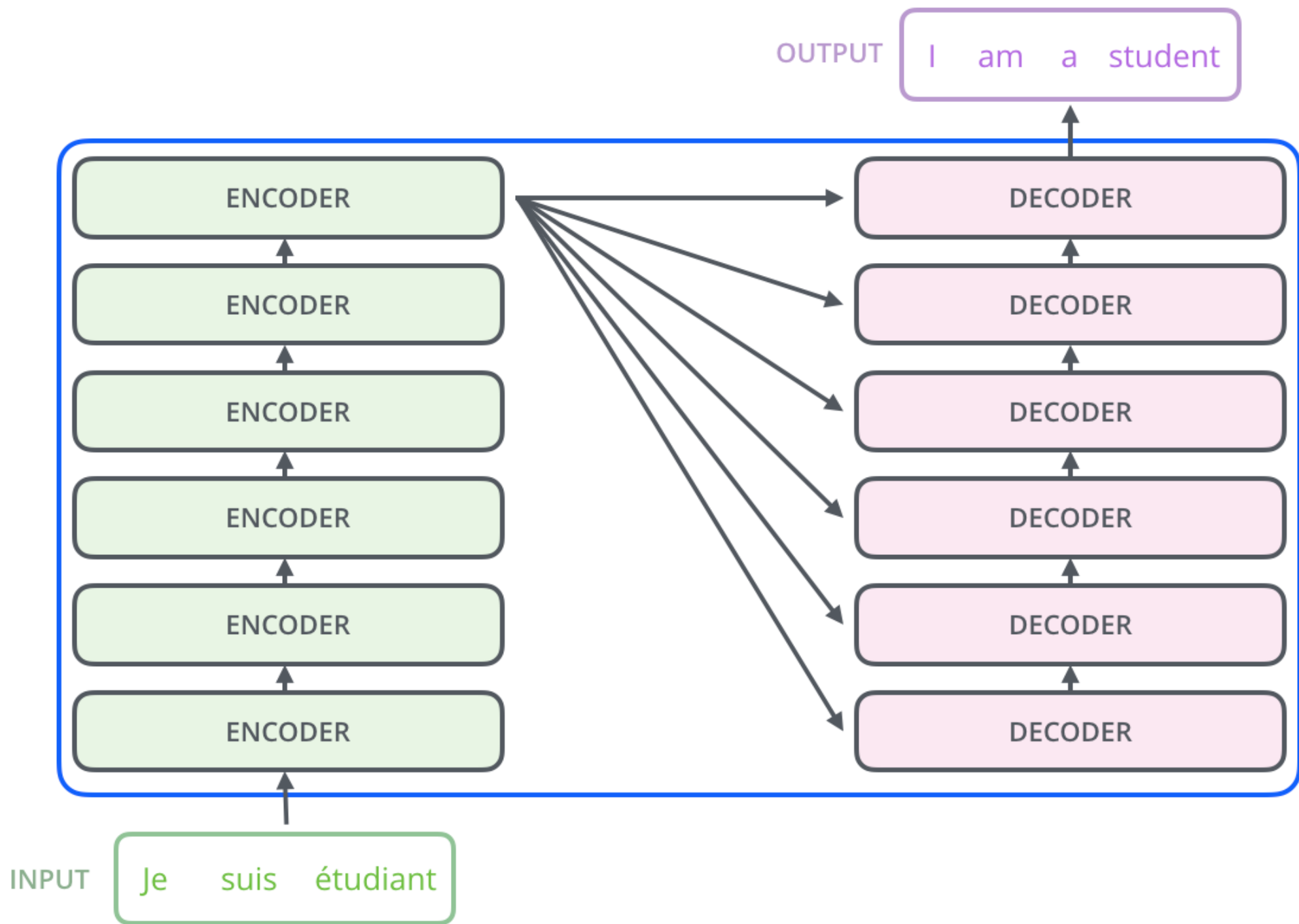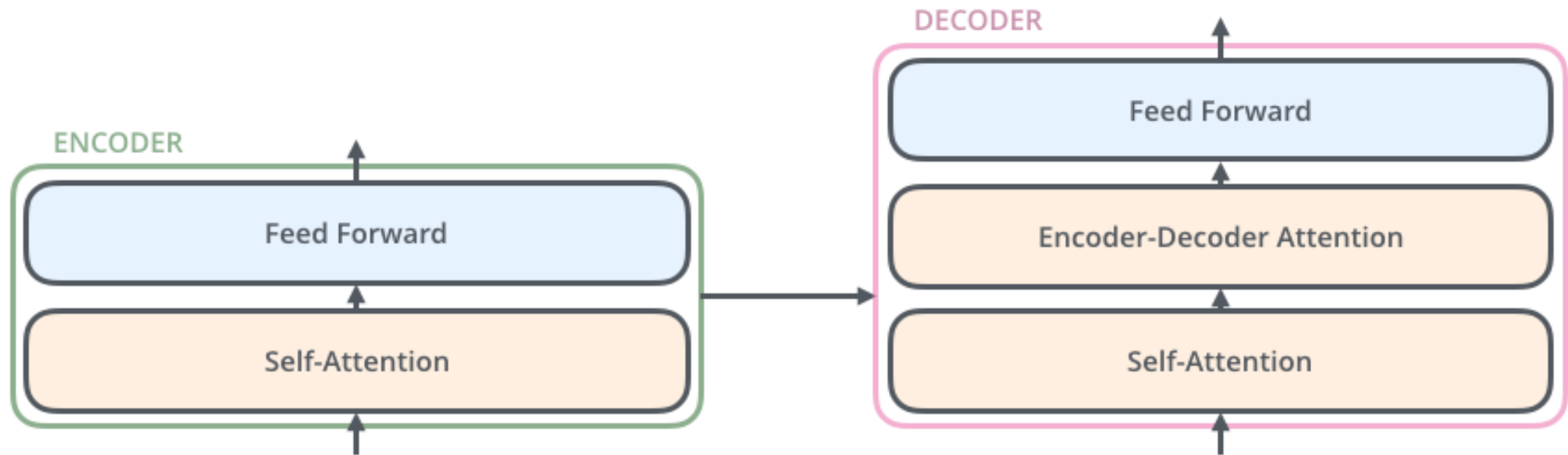- http://jalammar.github.io/illustrated-transformer/
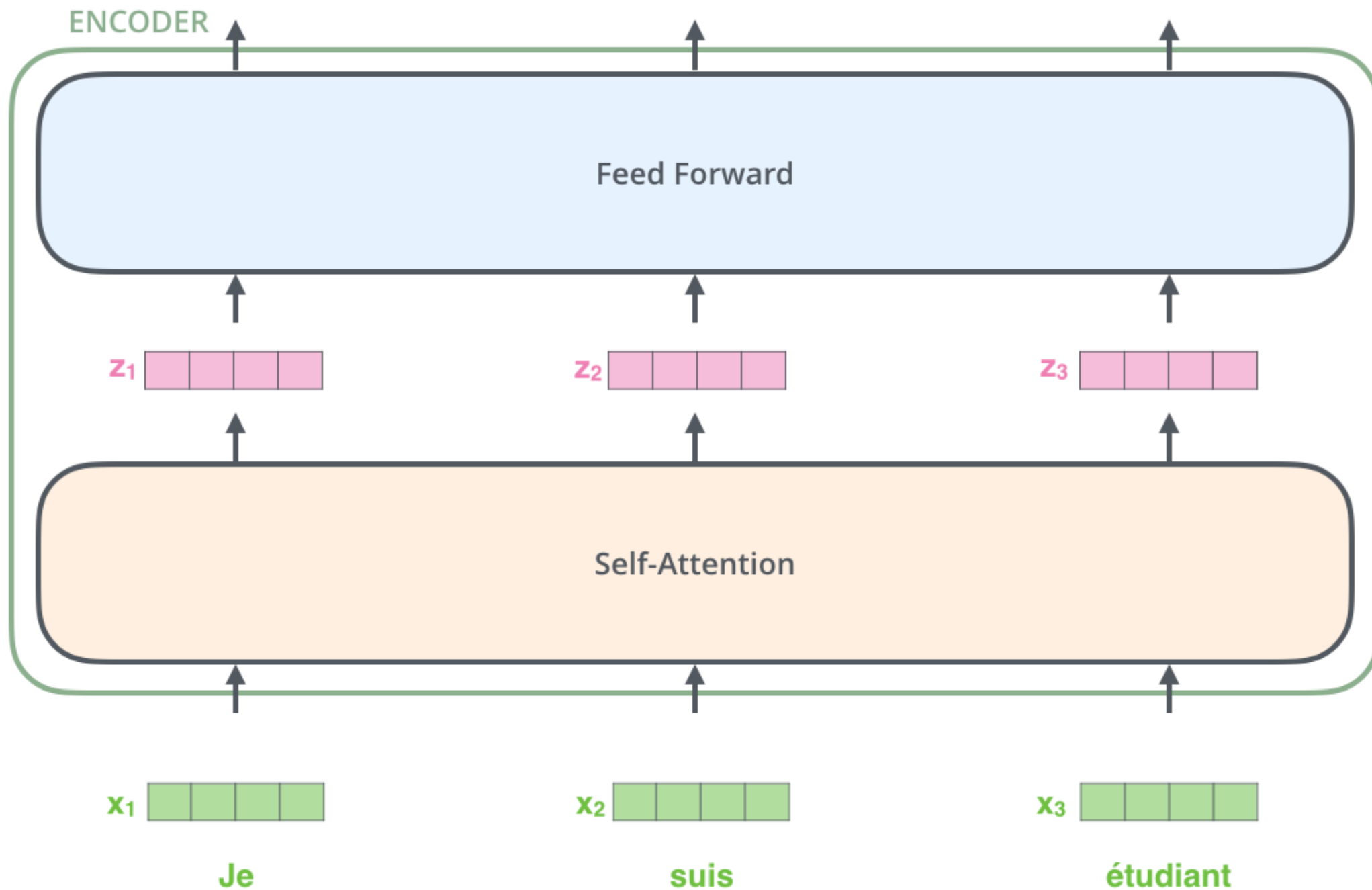
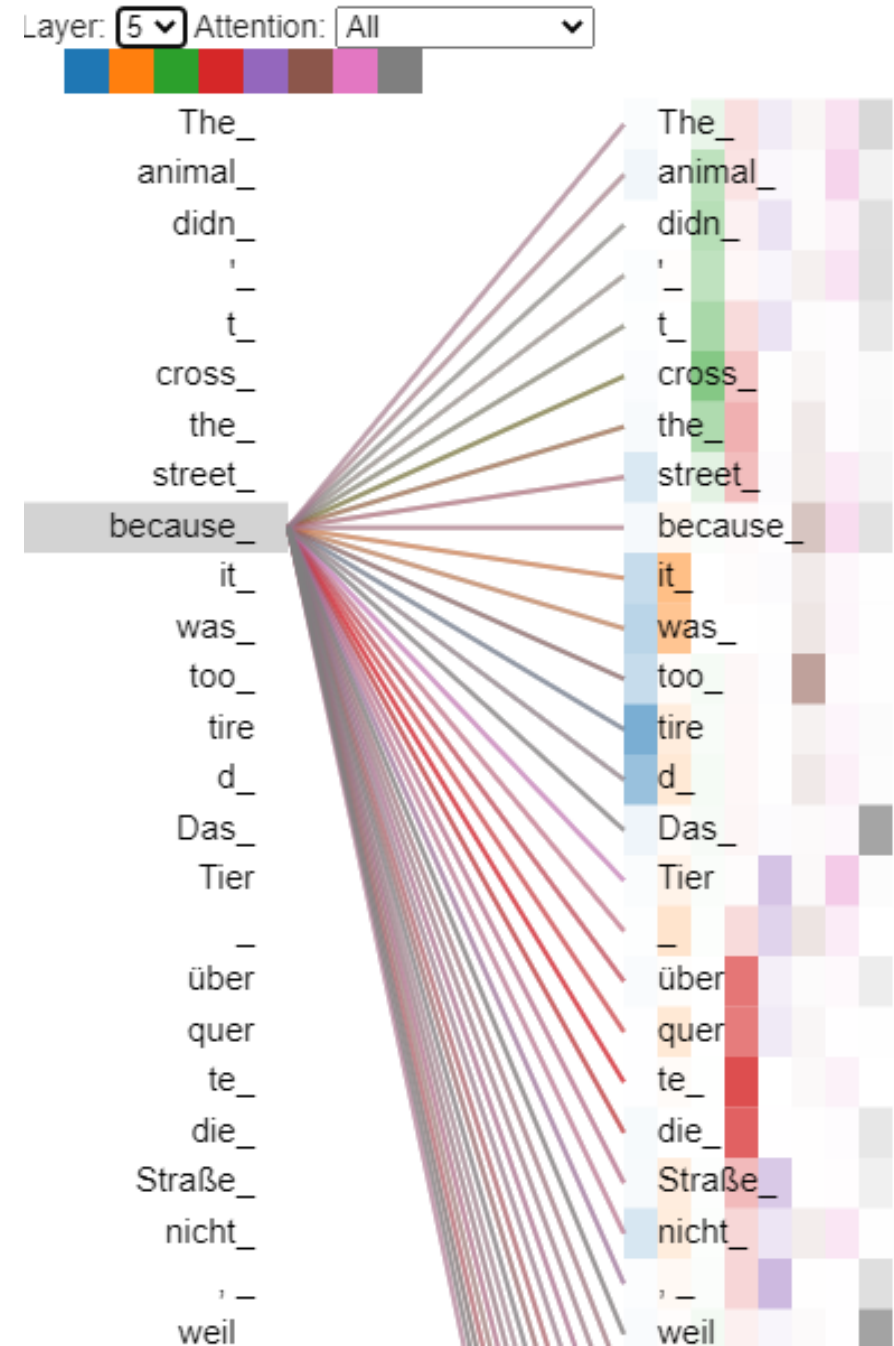# Encoder and Decoder

# Structure of the Encoder and Decoder

- Self-attention
- Encoder-decoder attention

# Tensor2Tensor Notebook

- https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello_t2t.ipynb

# Self-attention (query, key, value)
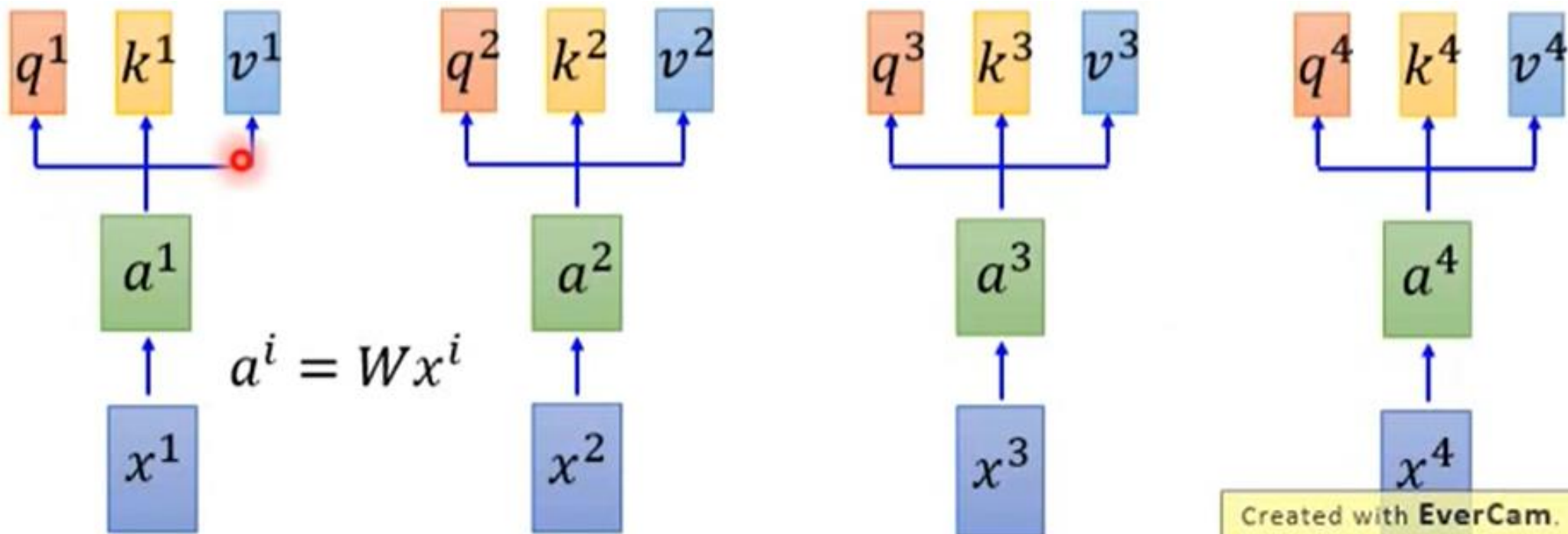
$q$: query (to match others)
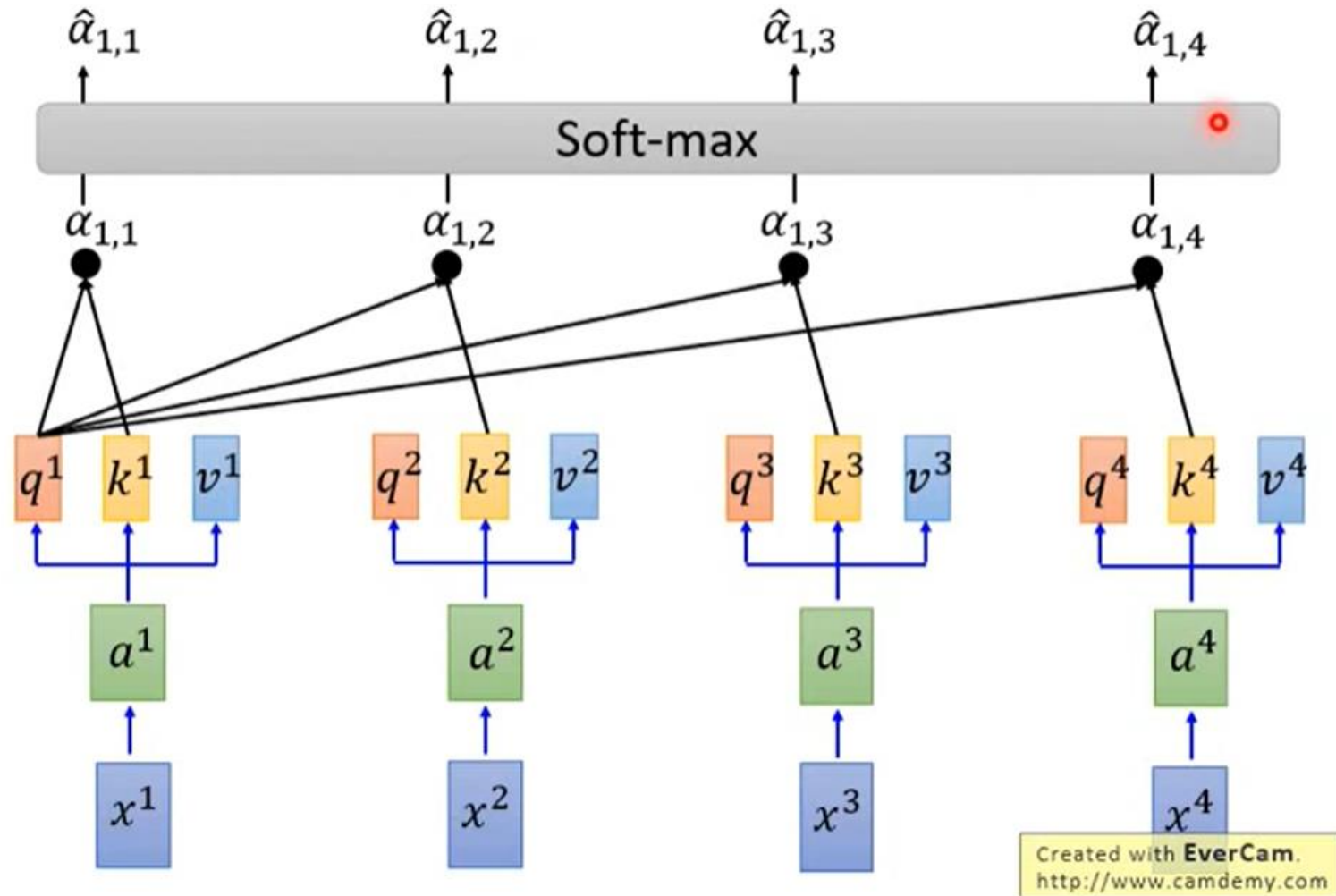$$q^i = W^q a^i$$

$k$: key (to be matched)
$$k^i = W^k a^i$$

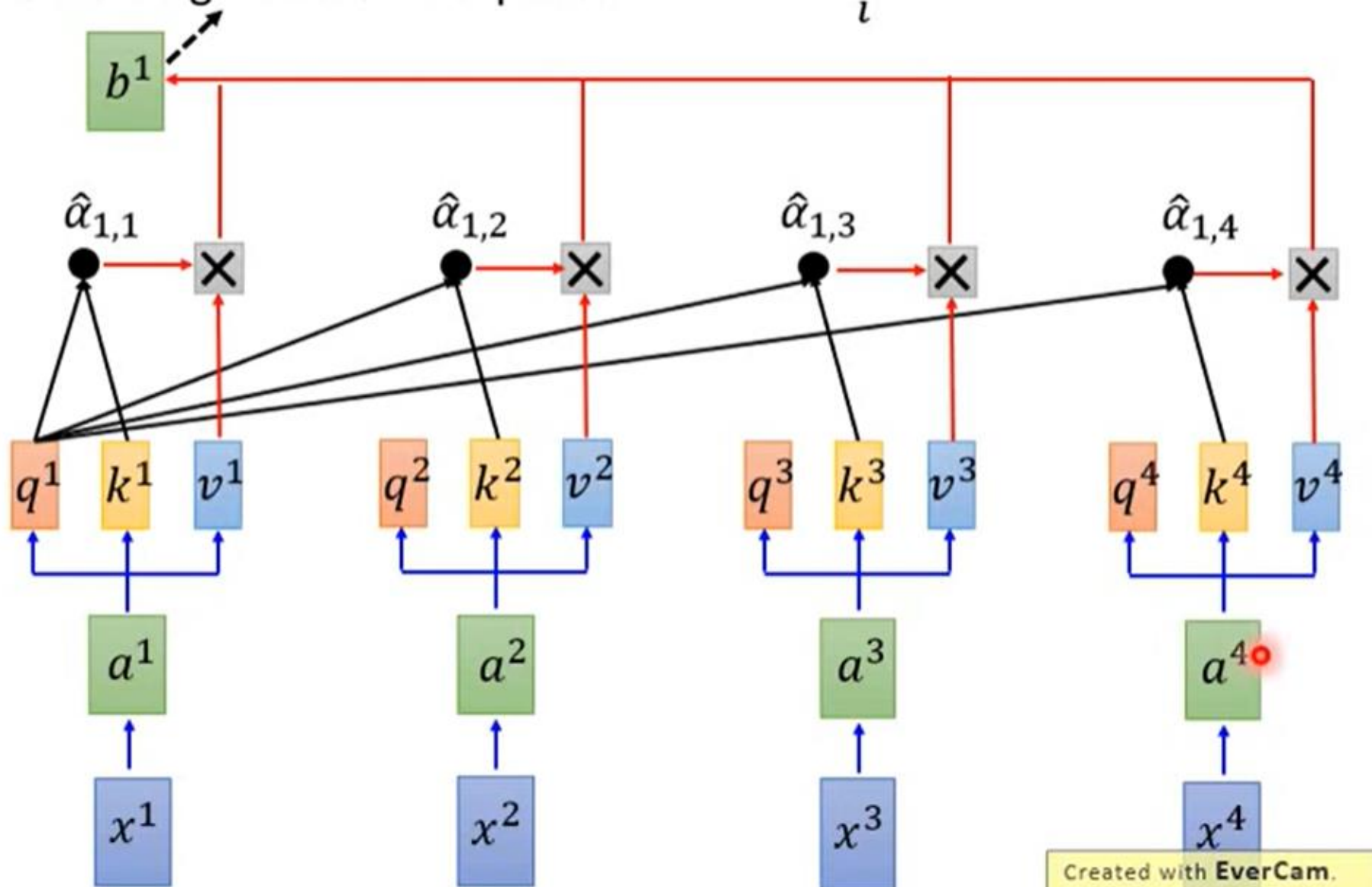$v$: information to be extracted
$$v^i = W^v a^i$$



$$a^i = W x^i$$

https://www.youtube.com/watch?v=ugWDIIOHtPA&t=1089s

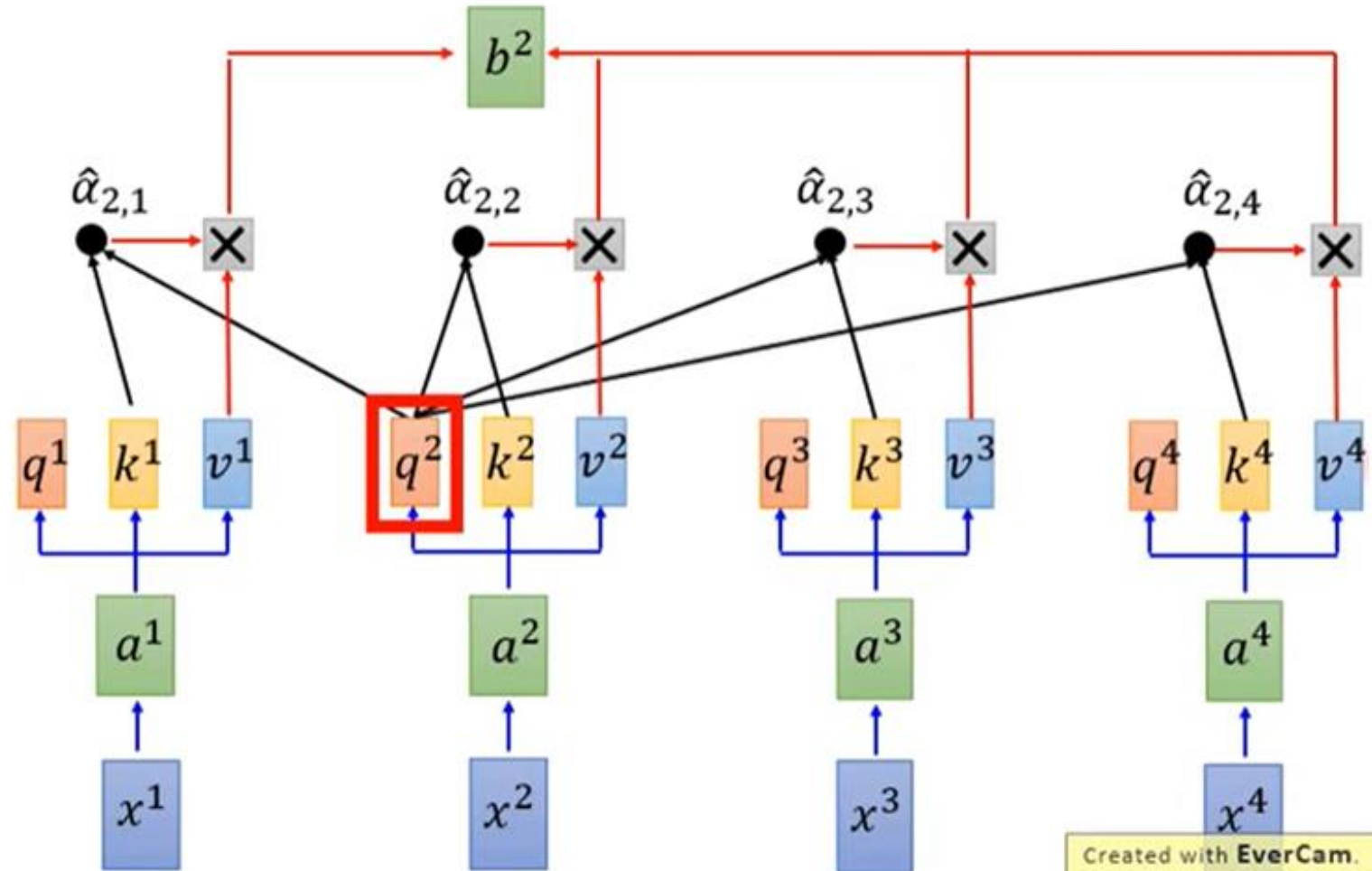# Self-attention

13

# Self-attention

Considering the whole sequence

$$b^1 = \sum_i \hat{\alpha}_{1,i} v^i$$

14

# Calculating $b^2$

$$b^2 = \sum_i \hat{\alpha}_{2,i} v^i$$

# Matrix Mutiplication
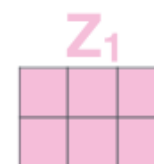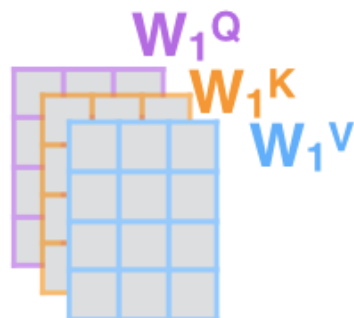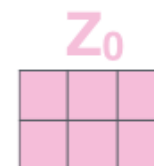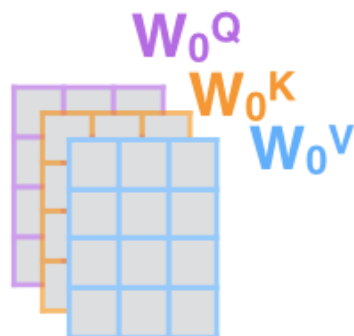
1) This is our input sentence*

2) We embed each word*

3) Split into 8 heads. We multiply $X$ or $R$ with weight matrices

4) Calculate attention using the resulting $Q$/$K$/$V$ matrices

5) Concatenate the resulting $Z$ matrices, then multiply with weight matrix $W^O$ to produce the output of the layer
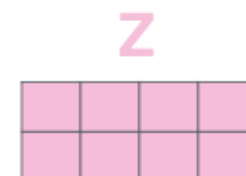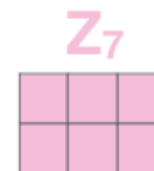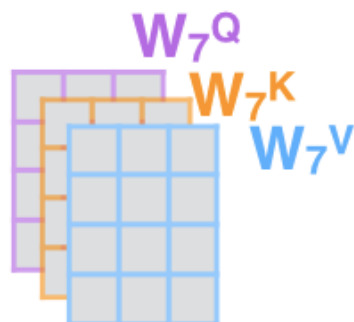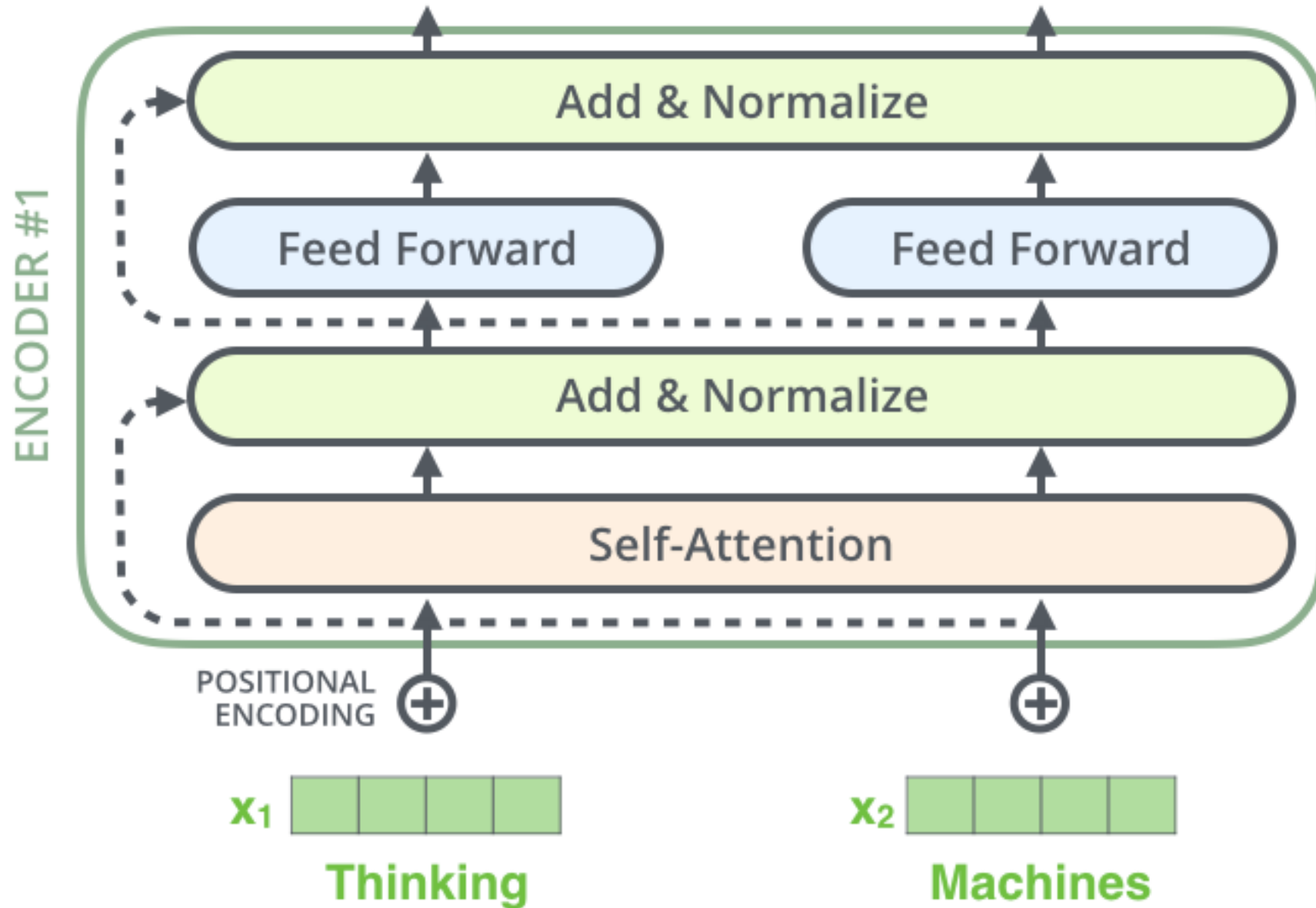
Thinking Machines

$X$

$W_0^Q$
$W_0^K$
$W_0^V$

$Q_0$
$K_0$
$V_0$

$Z_0$

$W^O$

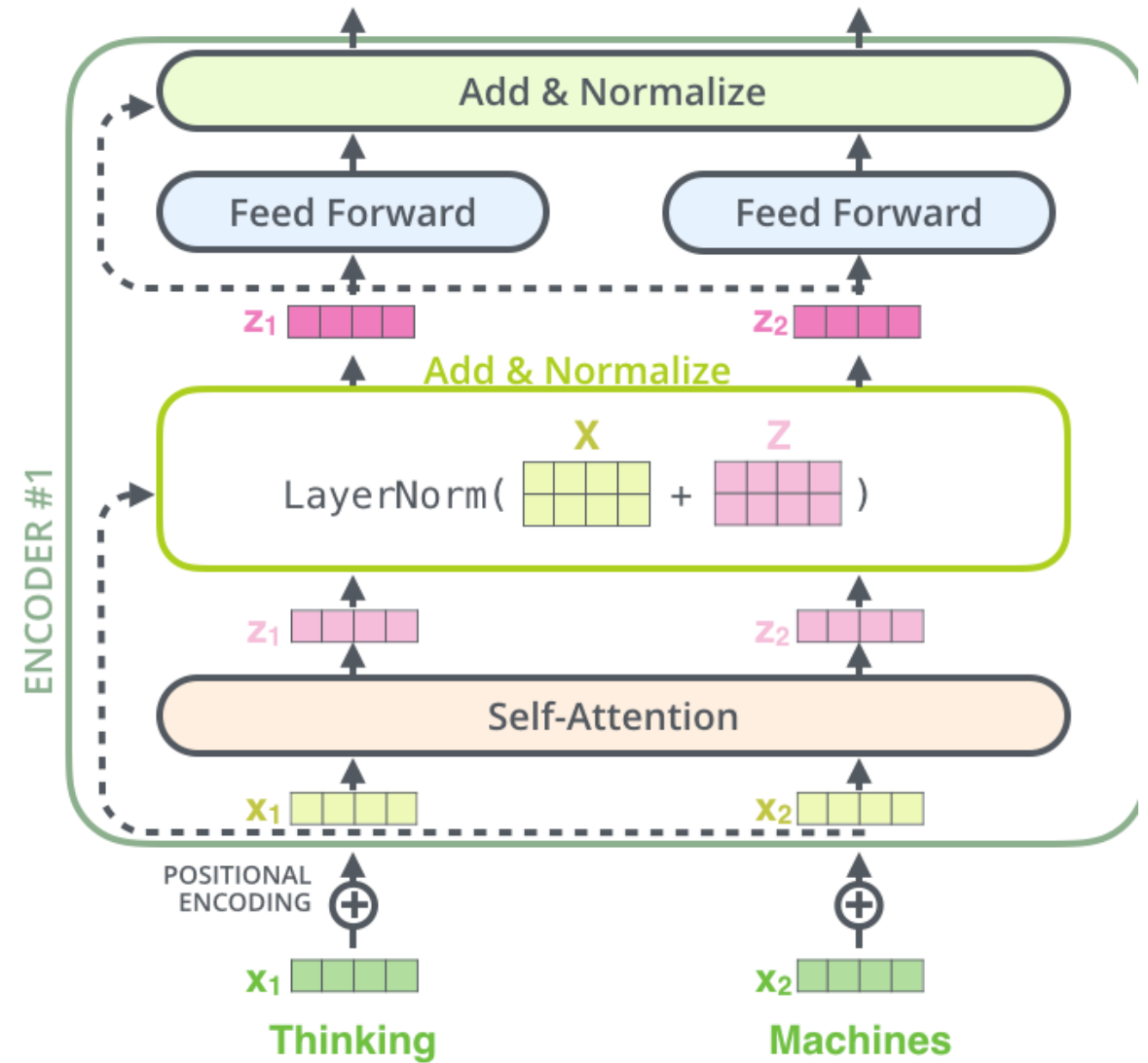* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one

$W_1^Q$
$W_1^K$
$W_1^V$

$Q_1$
$K_1$
$V_1$

$Z_1$

$Z$

...

...

...

$R$

$W_7^Q$
$W_7^K$
$W_7^V$

$Q_7$
$K_7$
$V_7$

$Z_7$

# Adding Residual Connections

# Layer Normalization

ENCODER #2

Add & Normalize

Feed Forward    Feed Forward

Add & Normalize

Self-Attention

ENCODER #1

Add & Normalize

Feed Forward    Feed Forward

Add & Normalize

Self-Attention

POSITIONAL ENCODING

$x_1$ Thinking

$x_2$ Machines

Softmax

Linear

DECODER #2

DECODER #1

Add & Normalize

Feed Forward    Feed Forward

Add & Normalize

Encoder-Decoder Attention

Add & Normalize

Self-Attention

# References

1. https://lilianweng.github.io/lil-log/2018/06/24/attention-attention.html

2. http://jalammar.github.io/illustrated-transformer/

3. Hong-Yi Lee, Transformer, 2019
   https://www.youtube.com/watch?v=ugWDIIOHtPA