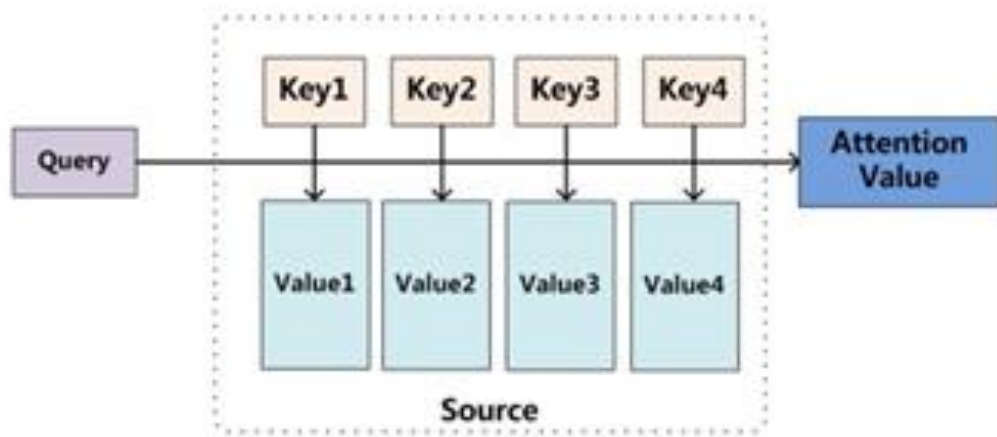# Attention is all you need

Steven Tang
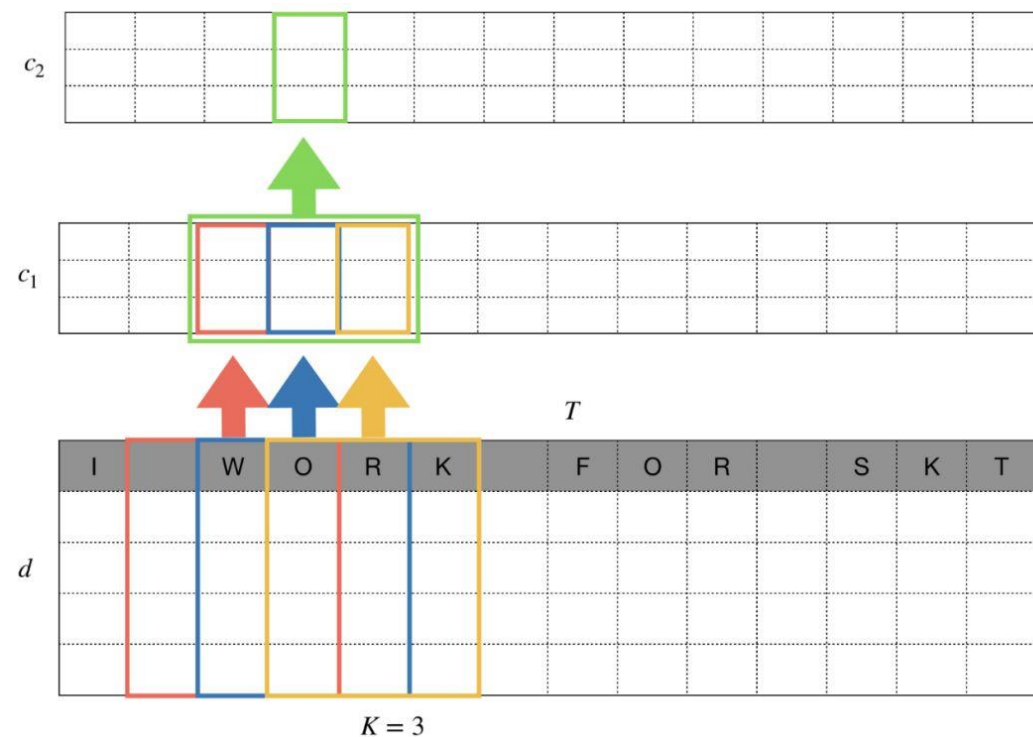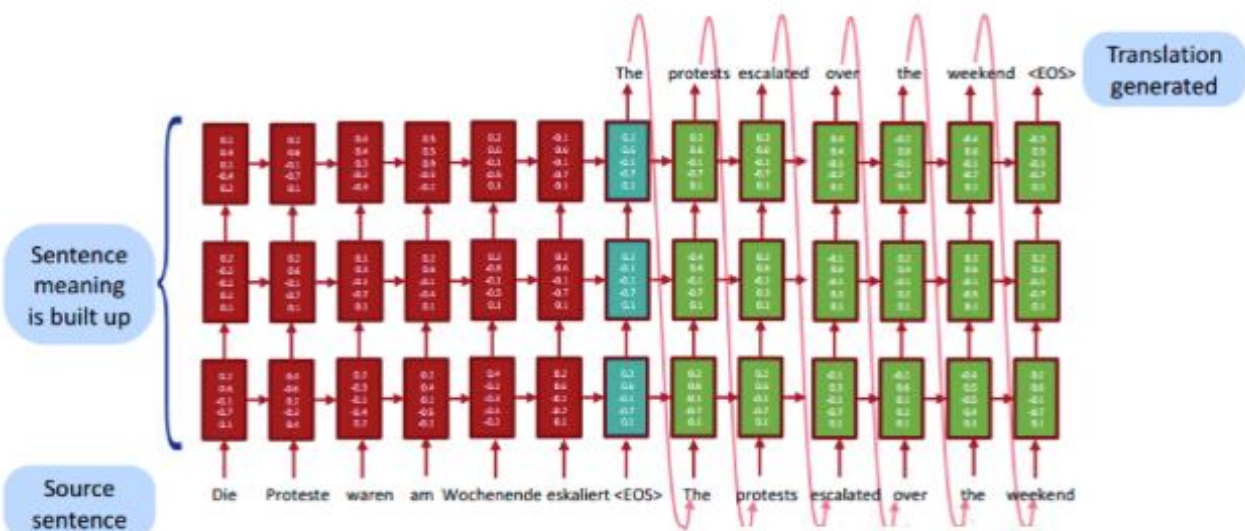
# 注意力机制回顾



- An attention function can be described as **mapping a query and a set of key-value pairs** to an output.
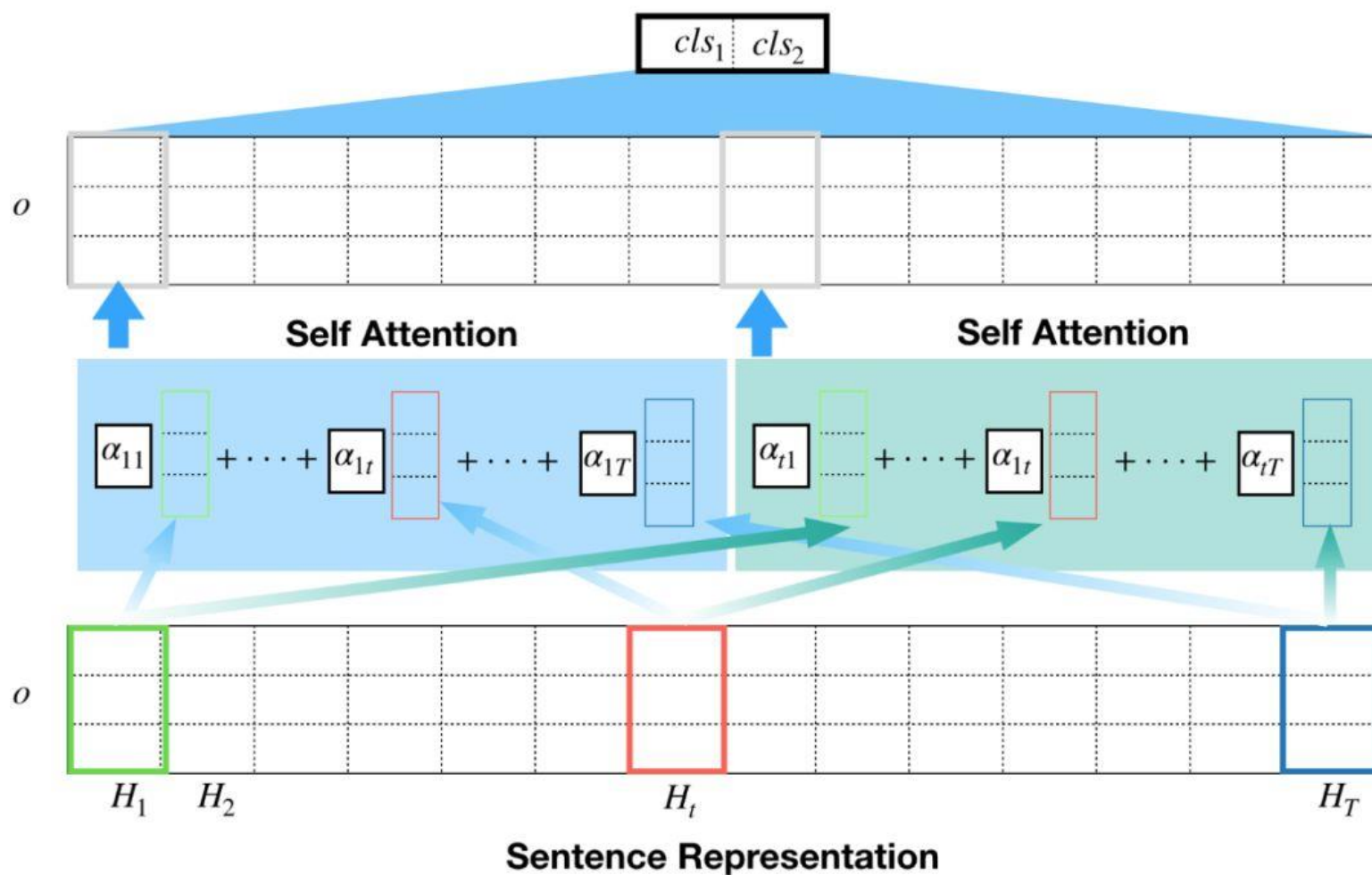
| Key1 | Key2 | Key3 | Key4 |

Query → Attention Value

Value1  Value2  Value3  Value4

Source

$$Attention(Query, Source) = \sum_{i=1}^{L_x} Similarity(Query, Key_i) \cdot Value_i$$
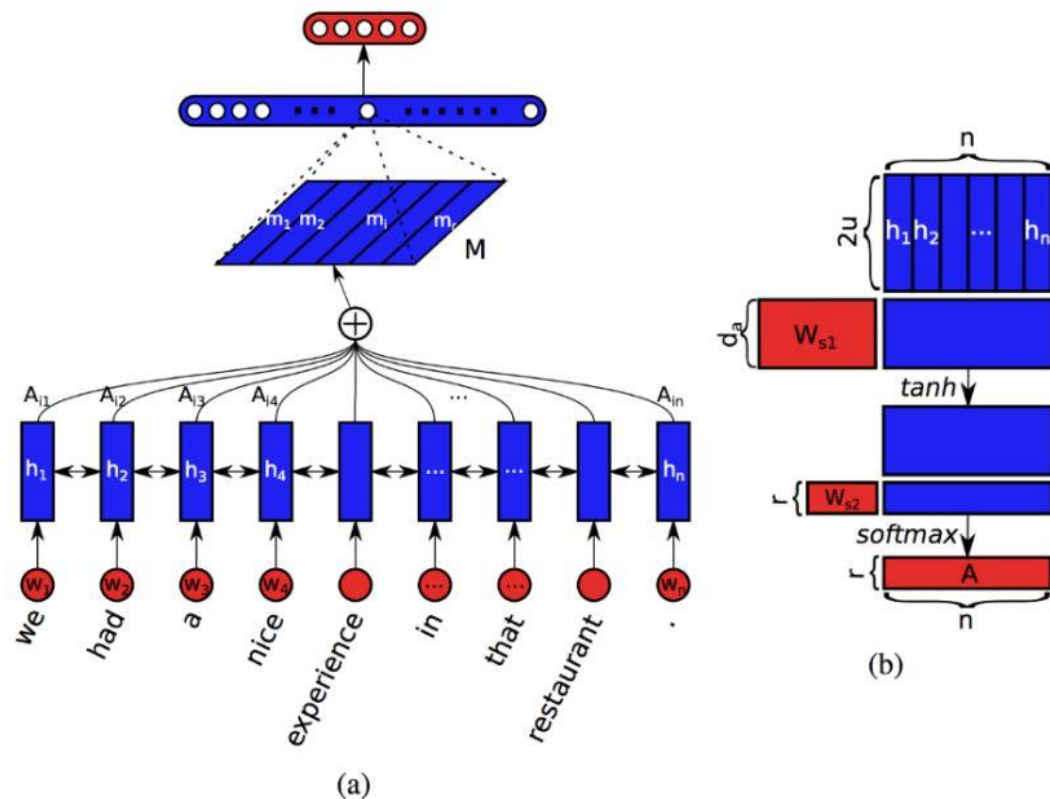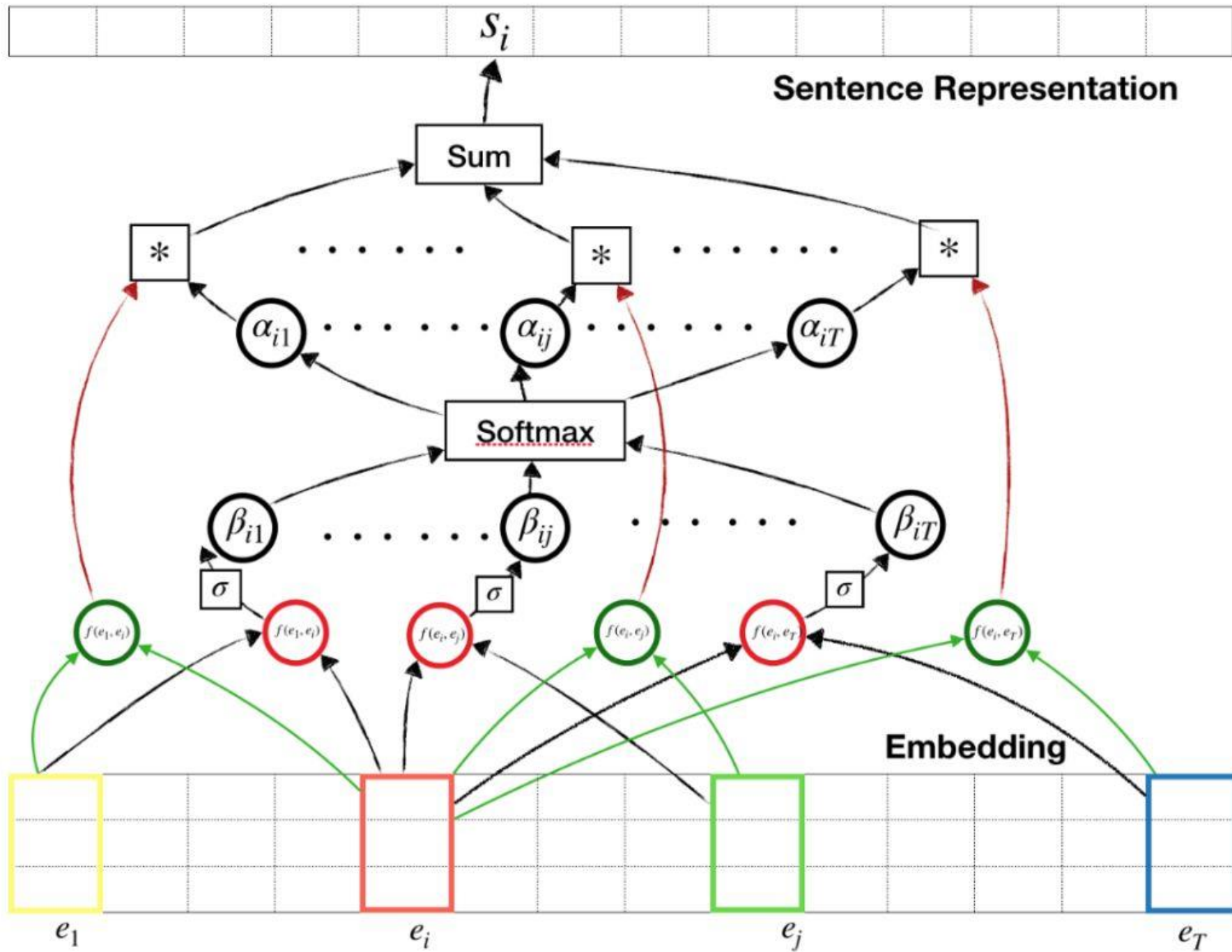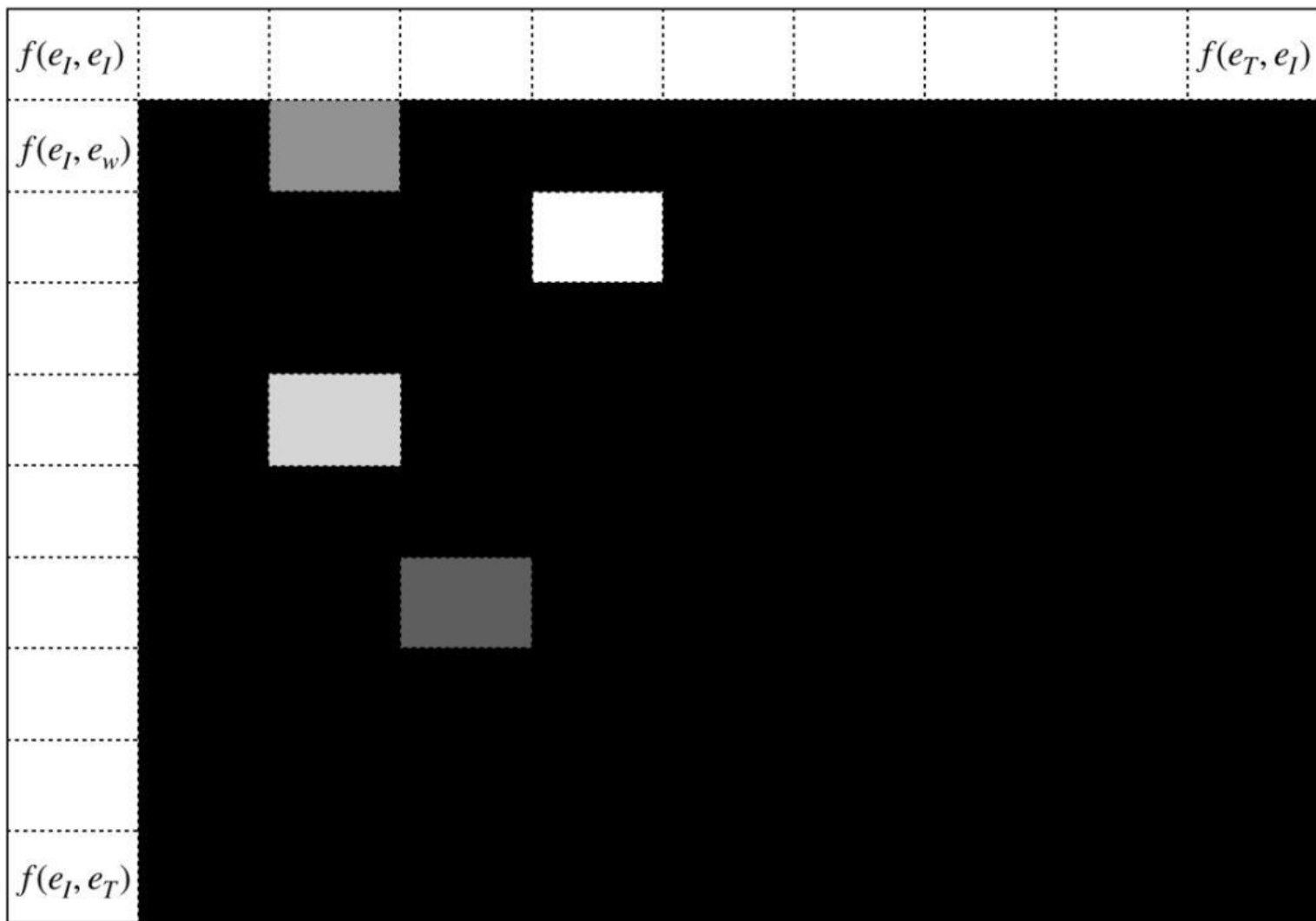
# CNN和RNN的问题

# 自注意力机制

# 自注意力橫跨出世



Figure 1: A sample model structure showing the sentence embedding model combined with a fully connected and softmax layer for sentiment analysis (a). The sentence embedding $M$ is computed as multiple weighted sums of hidden states from a bidirectional LSTM ($\mathbf{h_1}, ..., \mathbf{h_n}$), where the summation weights ($A_{i1}, ..., A_{in}$) are computed in a way illustrated in (b). Blue colored shapes stand for hidden representations, and red colored shapes stand for weights, annotations, or input/output.
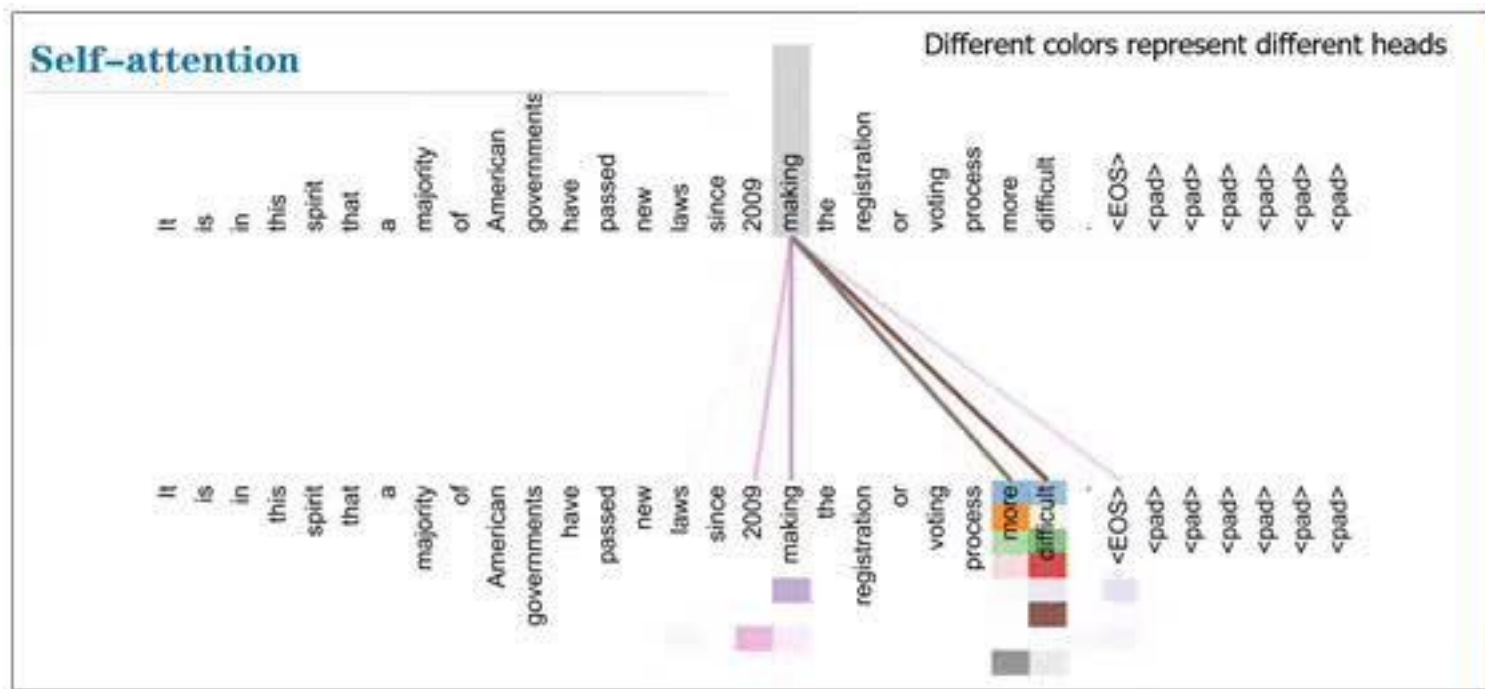
Sentence Representation

Embedding
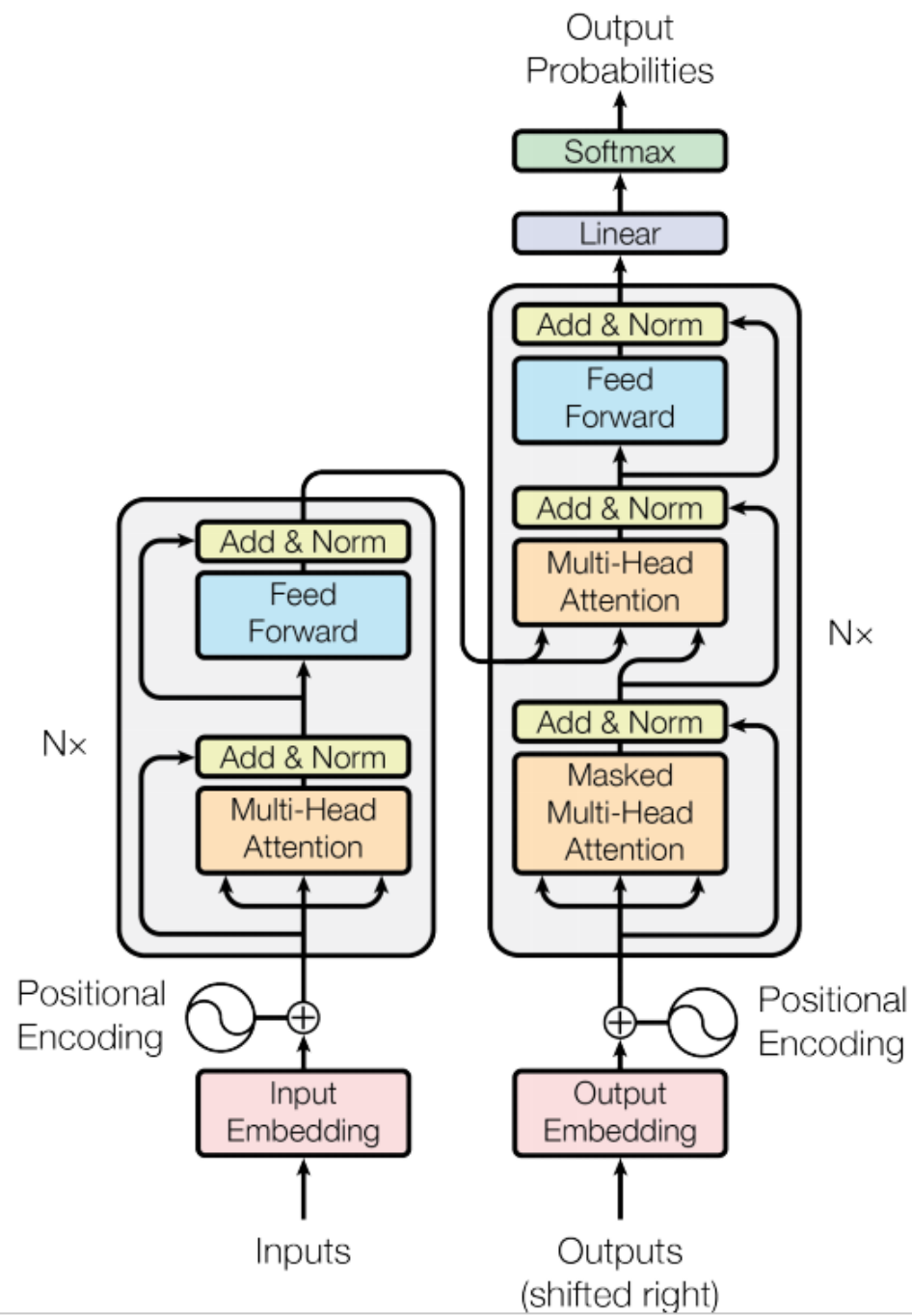
# 自注意力机制

# Self-attention（自注意力机制）



Fig. 6. The current word is in red and the size of the blue shade indicates the activation level. (Image source: Cheng et al., 2016)

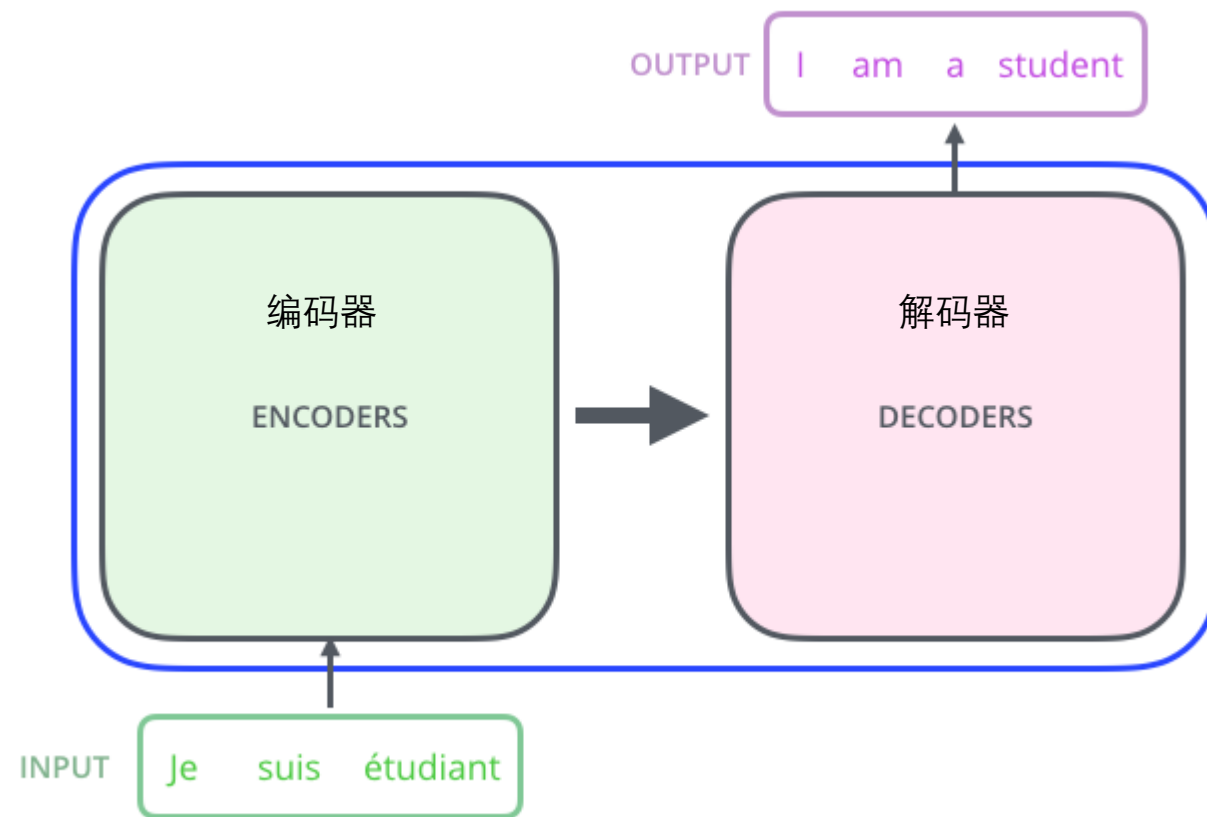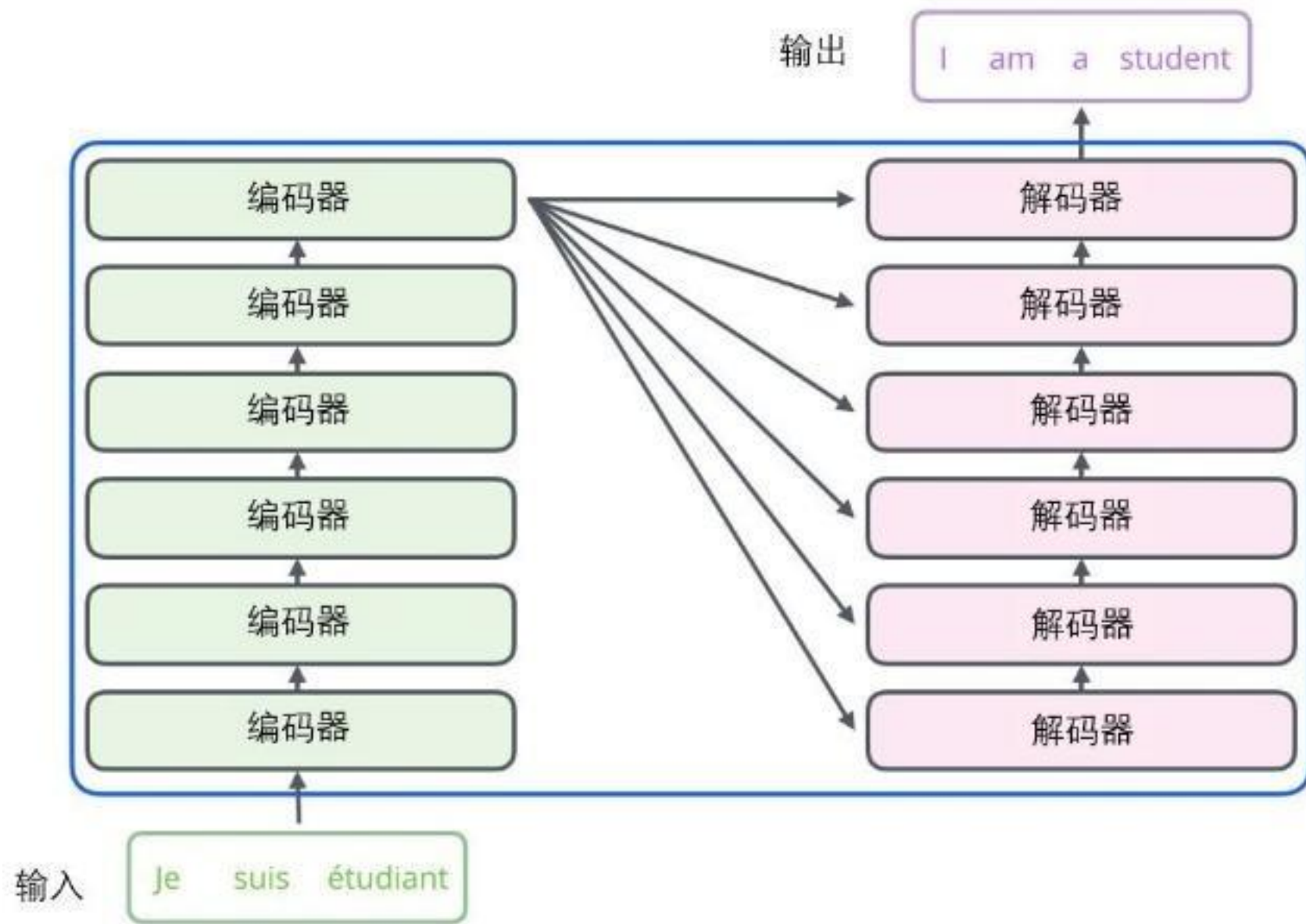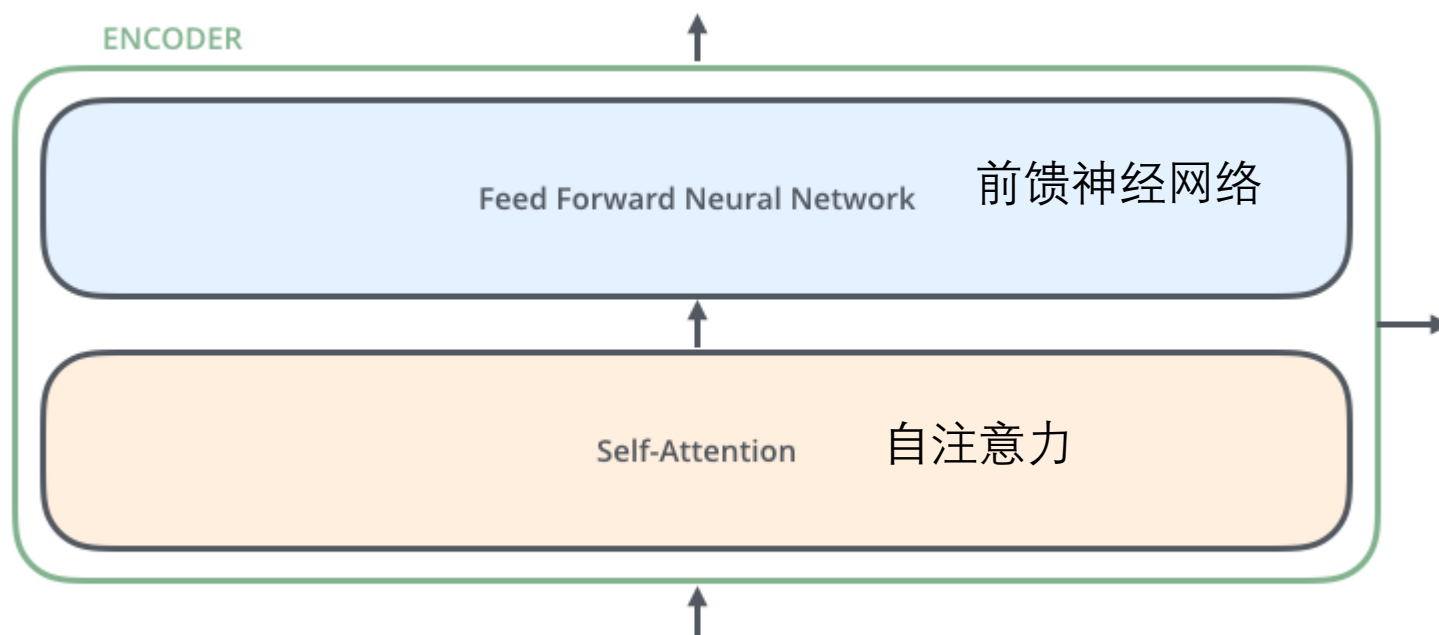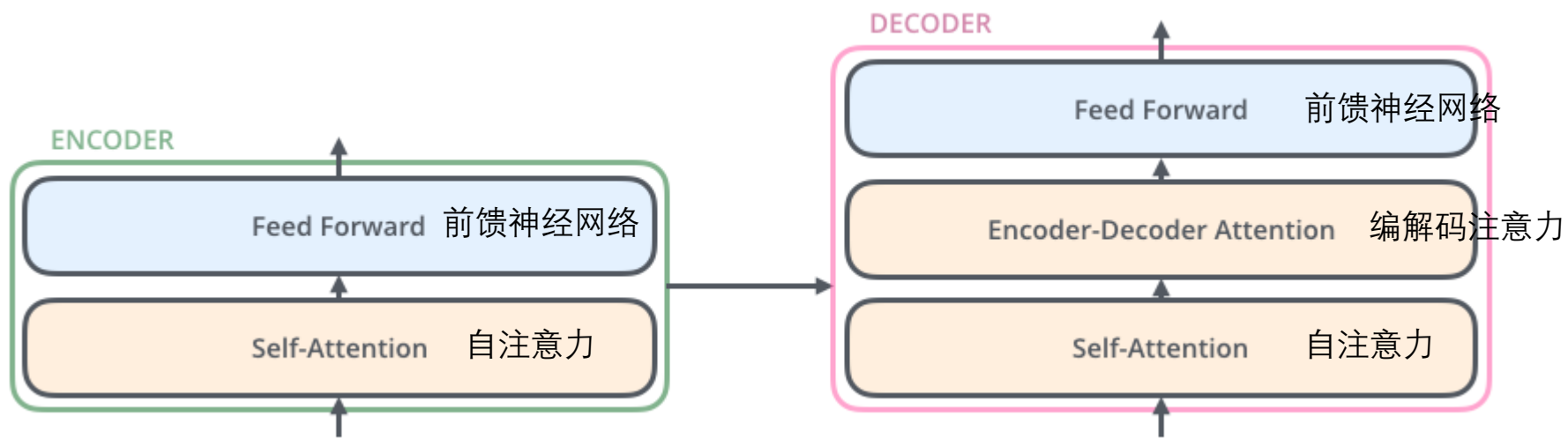# Self-attention（自注意力机制）

新的统
治者!

INPUT

Je suis étudiant

THE
TRANSFORMER

OUTPUT

I am a student

输出　I  am  a  student

编码器　　　　　　　→　解码器

编码器　　　　　　　　解码器

编码器　　　　　　　　解码器

编码器　　　　　　　　解码器

编码器　　　　　　　　解码器

编码器　　　　　　　　解码器

输入　Je  suis  étudiant

ENCODER

Feed Forward 前馈神经网络

Self-Attention 自注意力

DECODER

Feed Forward 前馈神经网络

Encoder-Decoder Attention 编解码注意力

Self-Attention 自注意力

ENCODER

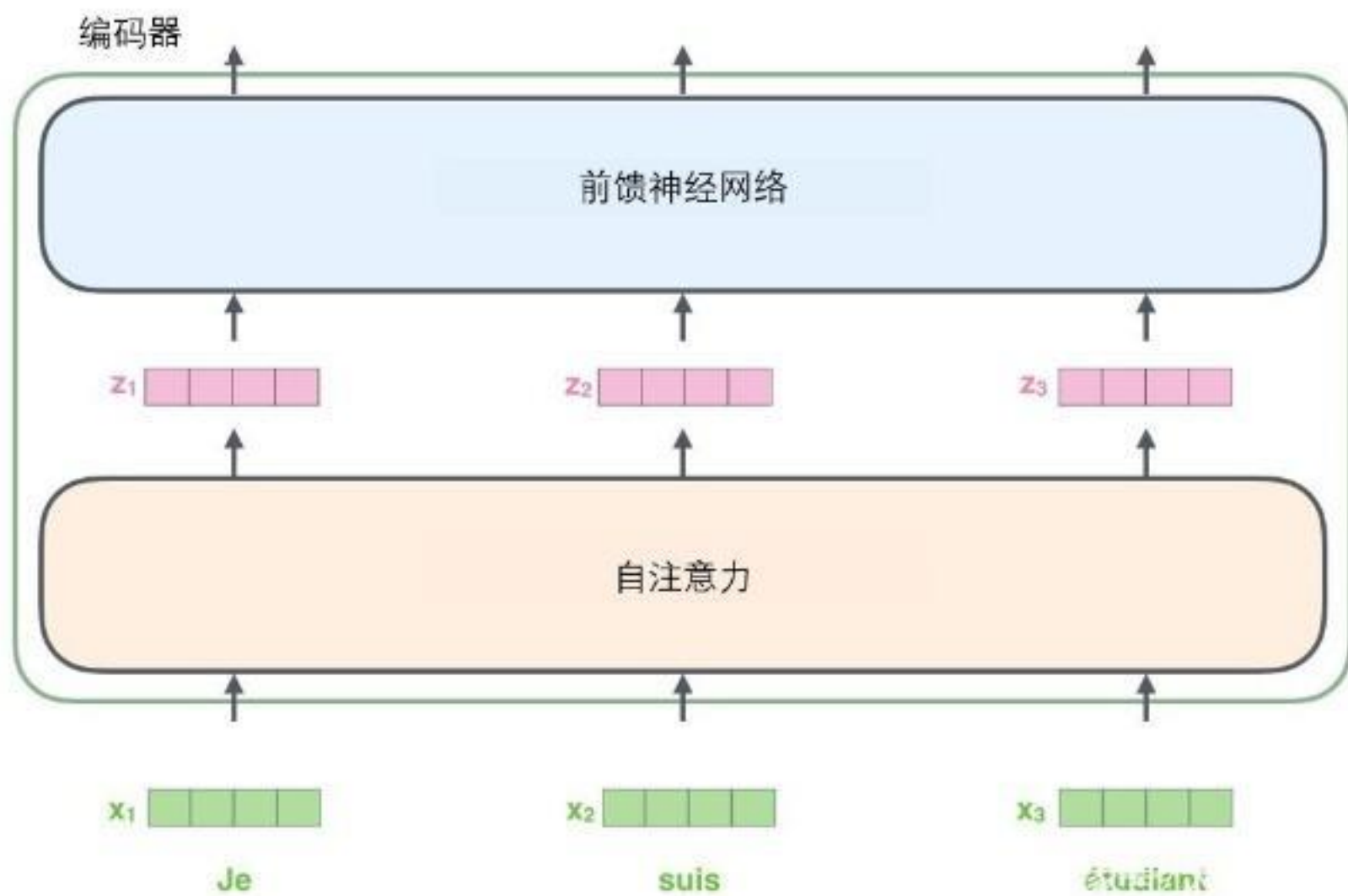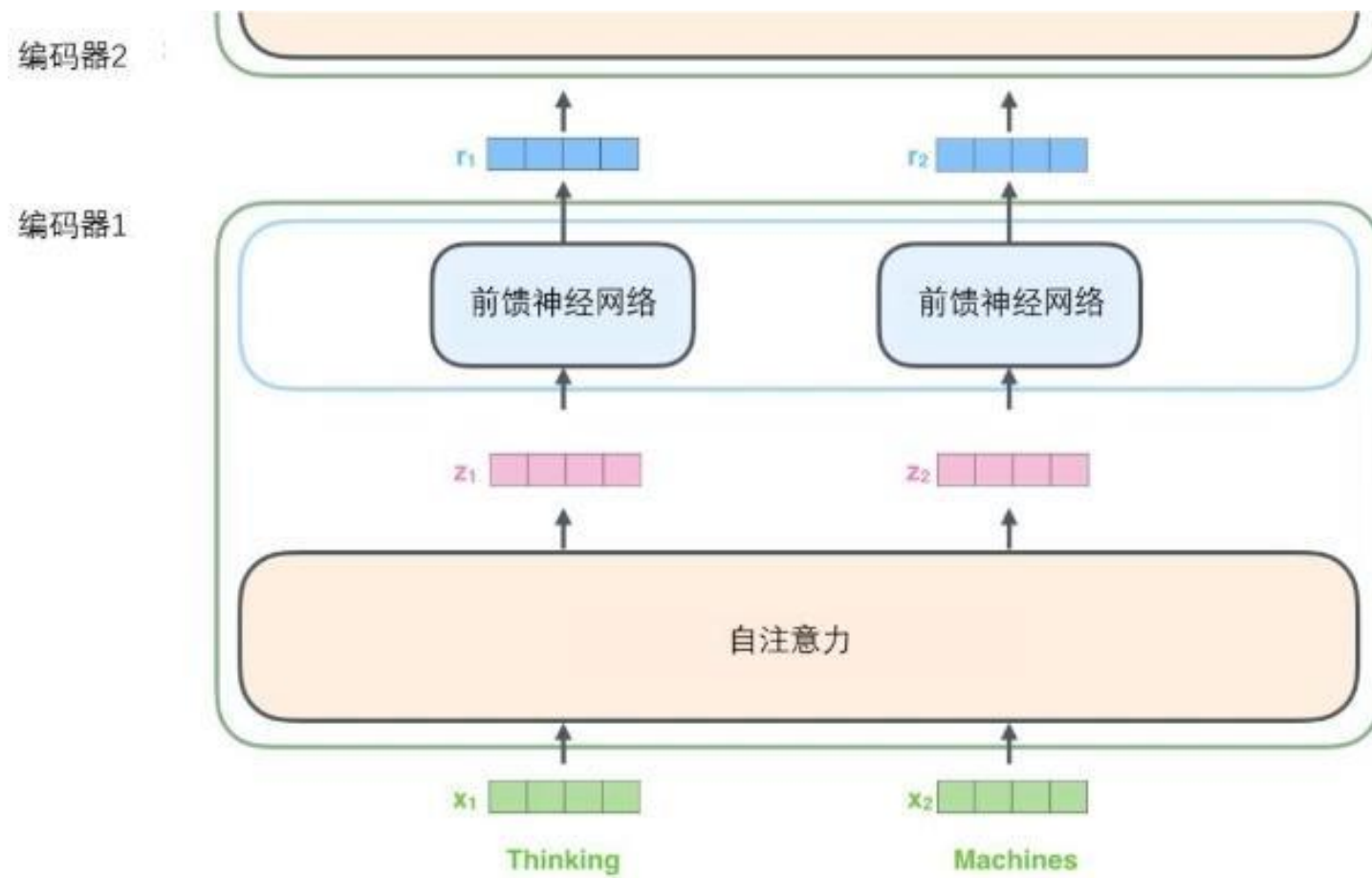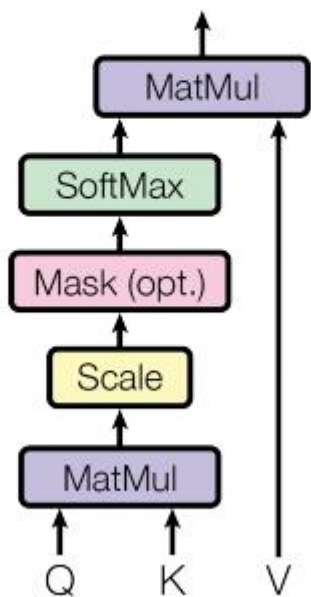Feed Forward Neural Network 前馈神经网络

Self-Attention 自注意力

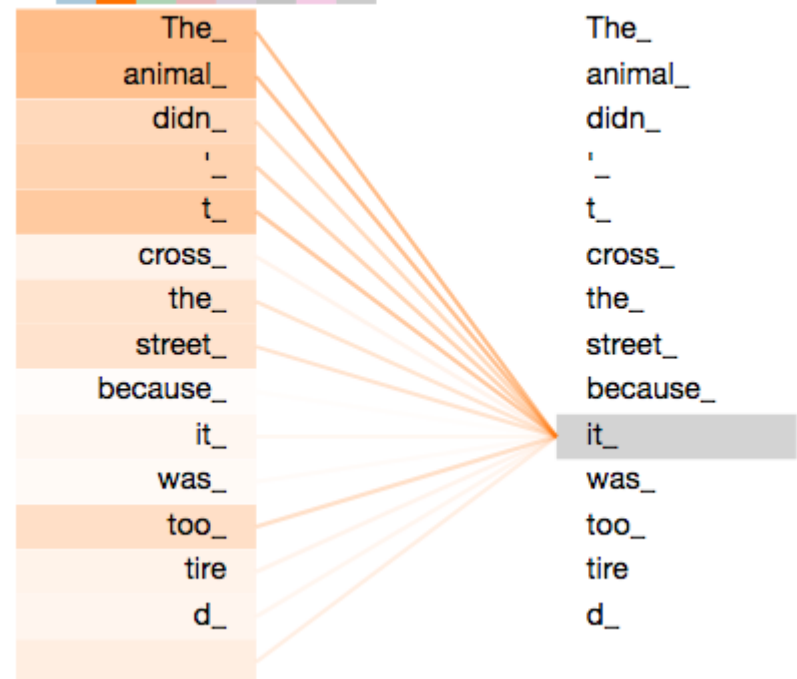# 降尺度的点乘注意力机制

**Scaled Dot-Product Attention**



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

Layer: 5 ▲▼ Attention: Input - Input ▼

| The_ | The_ |
| animal_ | animal_ |
| didn_ | didn_ |
| '_ | '_ |
| t_ | t_ |
| cross_ | cross_ |
| the_ | the_ |
| street_ | street_ |
| because_ | because_ |
| it_ | it_ |
| was_ | was_ |
| too_ | too_ |
| tire | tire |
| d_ | d_ |

Input    输入          **Thinking**          **Machines**

Embedding  嵌入    $X_1$ [ ][ ][ ][ ]    $X_2$ [ ][ ][ ][ ]

Queries    查询    $q_1$ [ ][ ][ ]    $q_2$ [ ][ ][ ]    $W^Q$

Keys    键    $k_1$ [ ][ ][ ]    $k_2$ [ ][ ][ ]    $W^K$

Values   值    $v_1$ [ ][ ][ ]    $v_2$ [ ][ ][ ]    $W^V$

| | Thinking | Machines |
|---|---|---|
| 输入 | | |
| 词嵌入 | $x_1$ | $x_2$ |
| 查询向量 | $q_1$ | $q_2$ |
| 键向量 | $k_1$ | $k_2$ |
| 值向量 | $v_1$ | $v_2$ |
| 打分 | $q_1 \cdot k_1 = 112$ | $q_1 \cdot k_2 = 96$ |

$$\text{softmax}\left( \frac{Q \times K^T}{\sqrt{d_k}} \right) V$$

$$= Z$$

| 输入 | | **Thinking** | | **Machines** | |
|---|---|---|---|---|---|
| 词嵌入 | $x_1$ | | $x_2$ | | |
| 查询向量 | $q_1$ | | $q_2$ | | |
| 键向量 | $k_1$ | | $k_2$ | | |
| 值向量 | $v_1$ | | $v_2$ | | |
| 打分 | | $q_1 \cdot k_1 = 112$ | | $q_1 \cdot k_2 = 96$ | |
| 除以8 （$\sqrt{d_k}$ ） | | 14 | | 12 | |
| Softmax | | 0.88 | | 0.12 | |
| softmax 乘以 值向量 | $v_1$ | | $v_2$ | | |
| 求和 | $z_1$ | | $z_2$ | | |

# 多头注意力机制



**Multi-Head Attention**

# 多头注意力机制

# 多头注意力机制

X

Thinking
Machines

分别计算这八个头的注意力值

Calculating attention separately in
eight different attention heads

ATTENTION
HEAD #0

ATTENTION
HEAD #1

...

ATTENTION
HEAD #7

$Z_0$

$Z_1$

$Z_7$

# 多头注意力机制

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, ..., \text{head}_h)W^O$$

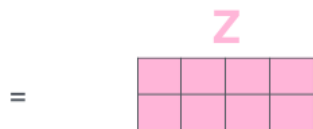$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

与W0相乘，这个矩阵也会和模型一起训练

将所有的注意力头拼接

1) Concatenate all the attention heads

$Z_0$　$Z_1$　$Z_2$　$Z_3$　$Z_4$　$Z_5$　$Z_6$　$Z_7$

2) Multiply with a weight matrix $W^O$ that was trained jointly with the model

X

最后的Z矩阵会捕捉所有注意力头的特征，然后讲这个Z送往前馈网络。

3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN

Z

=

$W^O$

1) This is our input sentence*

输入序列

Thinking Machines

2) We embed each word*

嵌入矩阵

**X**

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

使用8个注意力头

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

4) Calculate attention using the resulting Q/K/V matrices
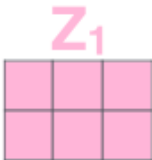
利用QKV矩阵计算注意力值

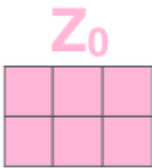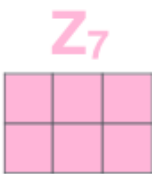$Q_0$
$K_0$
$V_0$

5) Concatenate the resulting Z matrices, then multiply with weight matrix $W^O$ to produce the output of the layer
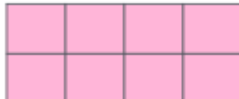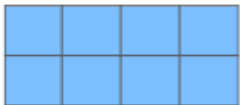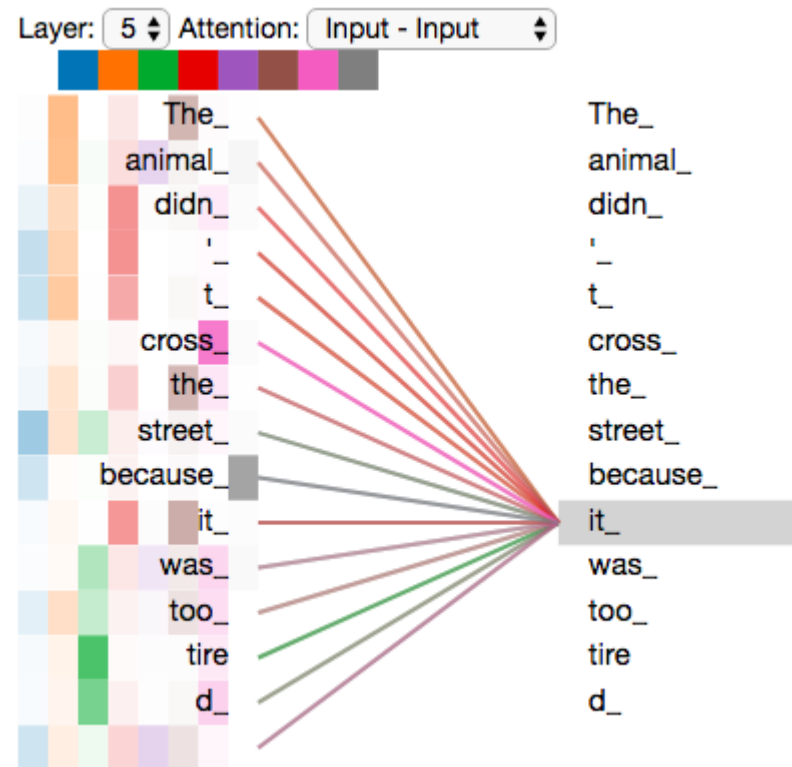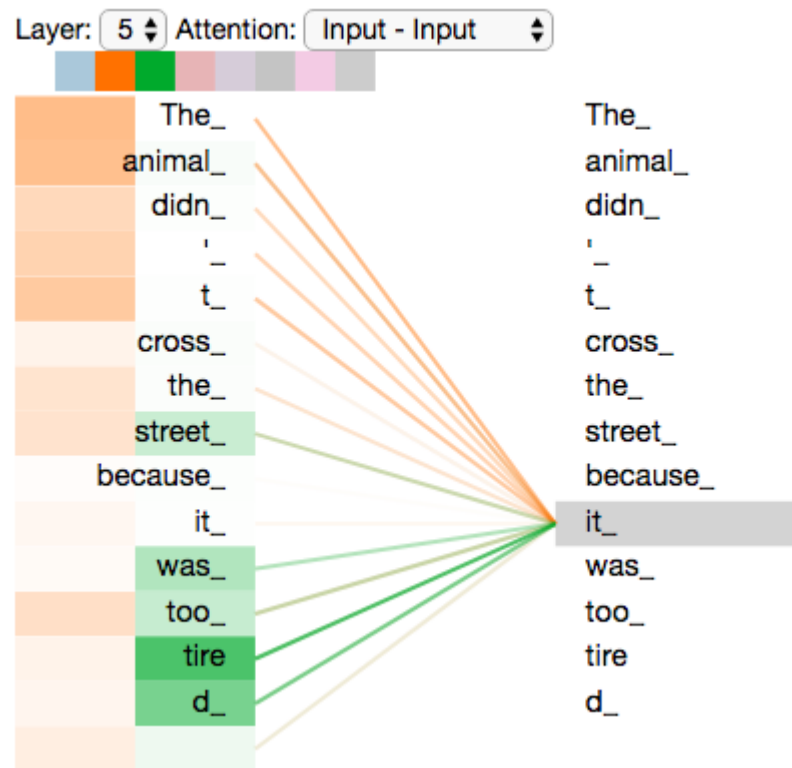
将拼接后的矩阵和W0相乘得到Z

$Z_0$

**$W^O$**

* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one
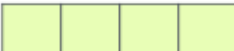
对处理0号注意力头的所有编码器，我们不需要嵌入，直接从0号下面的编码器的输出开始

**R**

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

$Q_1$
$K_1$
$V_1$

$Z_1$

**Z**

...

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

...

$Q_7$
$K_7$
$V_7$

...

$Z_7$

| 0 | 0 | 1 | 1 |

| 0.84 | 0.0001 | 0.54 | 1 |

| 0.91 | 0.0002 | -0.42 | 1 |

+

+

+

EMBEDDINGS
嵌入向量

$X_1$

$X_2$

$X_3$

输入序列
INPUT

Je

suis

étudiant

$$PE_{(pos,2i)} = sin(pos/10000^{2i/d_{model}})$$
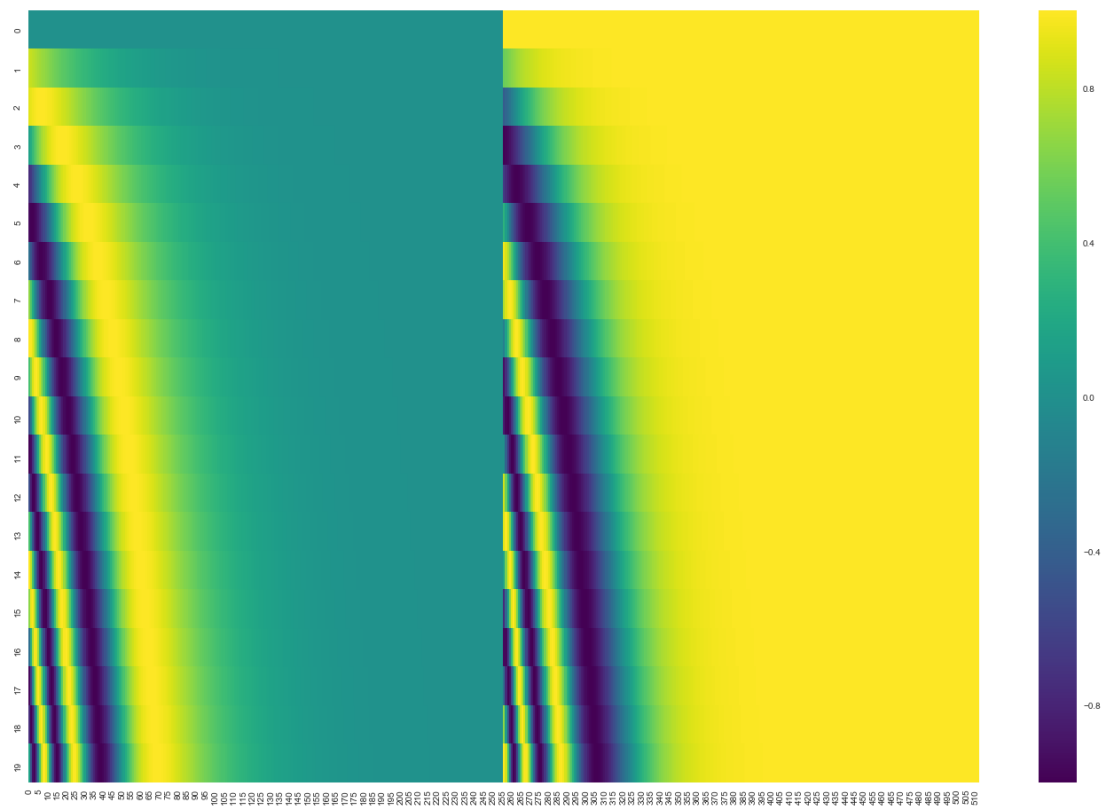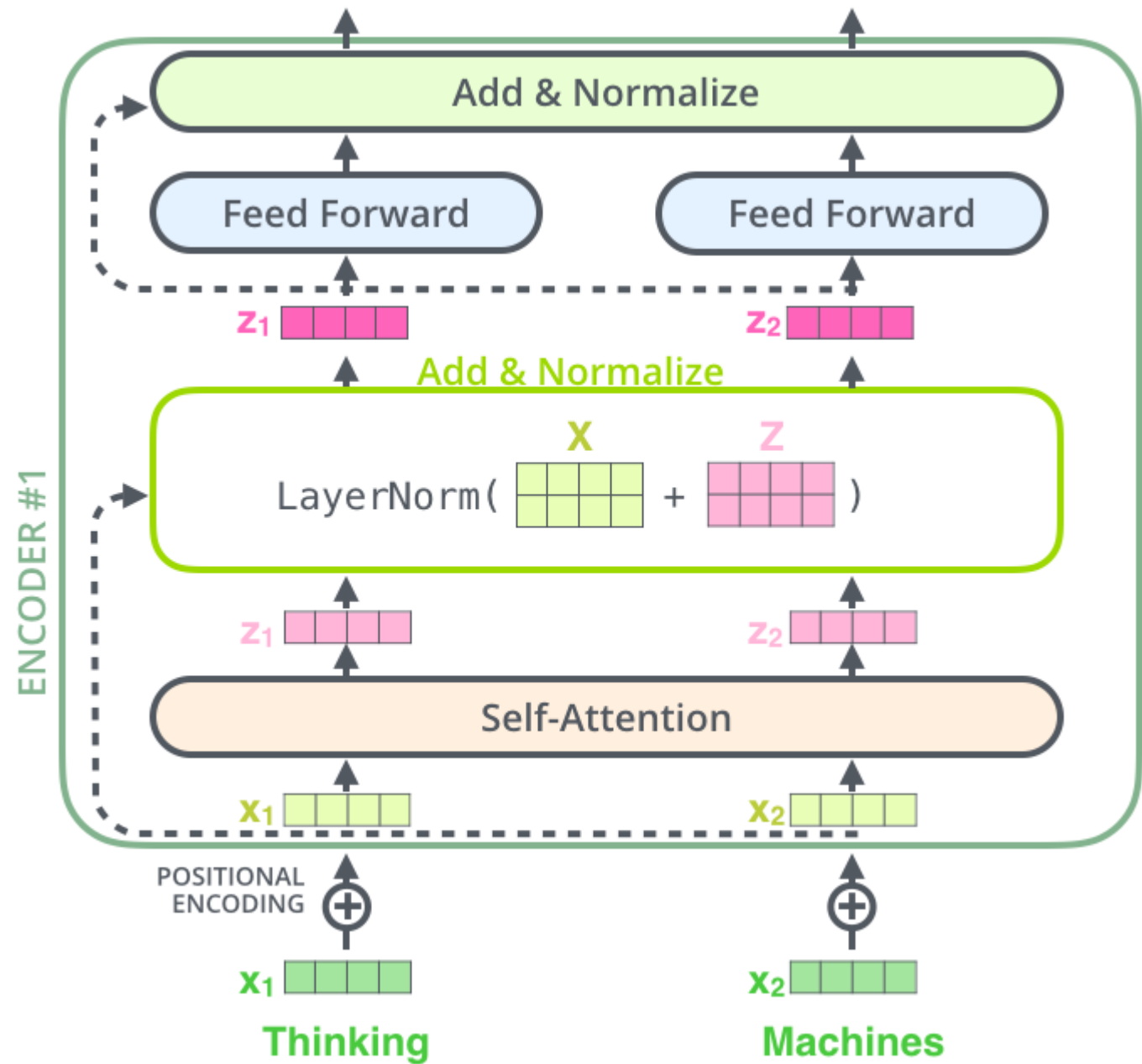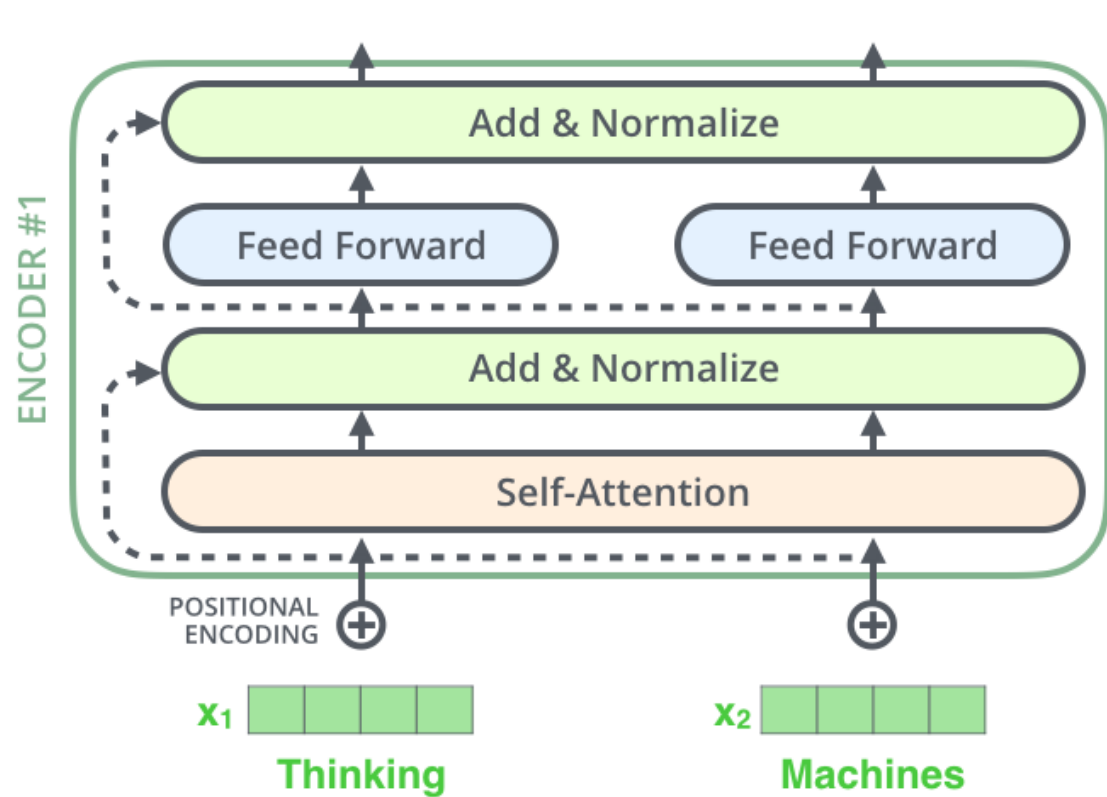
$$PE_{(pos,2i+1)} = cos(pos/10000^{2i/d_{model}})$$

**ENCODER #1** (left diagram)

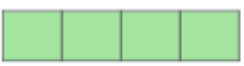Add & Normalize

Feed Forward | Feed Forward

Add & Normalize

Self-Attention

POSITIONAL ENCODING ⊕    ⊕

$x_1$ Thinking    $x_2$ Machines

**ENCODER #1** (right diagram)

Add & Normalize

Feed Forward | Feed Forward

$z_1$    $z_2$

Add & Normalize

$$\texttt{LayerNorm}(\quad X \quad + \quad Z \quad)$$

$z_1$    $z_2$

Self-Attention

$x_1$    $x_2$

POSITIONAL ENCODING ⊕    ⊕

$x_1$ Thinking    $x_2$ Machines

Decoding time step: ⬤1 2 3 4 5 6          OUTPUT

Linear + Softmax

ENCODER                    DECODER

ENCODER                    DECODER
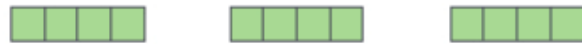
EMBEDDING
WITH TIME
SIGNAL

EMBEDDINGS

INPUT        Je        suis      étudiant

Which word in our vocabulary
is associated with this index?

am

Get the index of the cell
with the highest value
(`argmax`)

5

**log_probs**

0 1 2 3 4 5                    … vocab_size

Softmax

**logits**

0 1 2 3 4 5                    … vocab_size

Linear

Decoder stack output