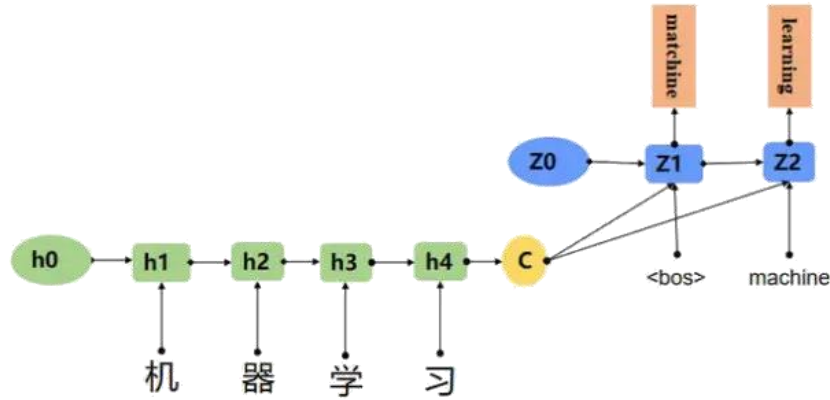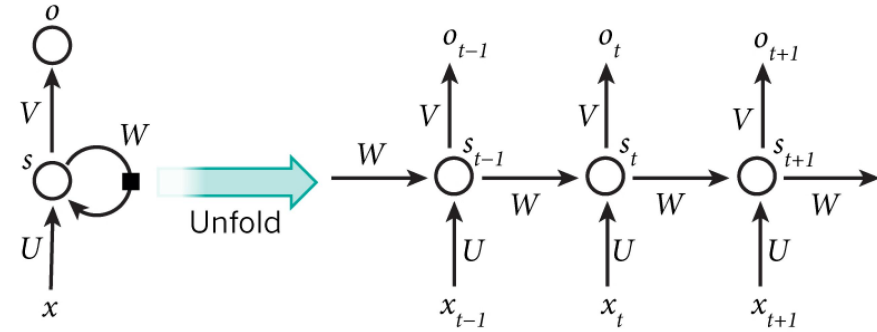# Attention Is All You Need 论文导读
# &
# Transformer详解

# Why Transformer ?
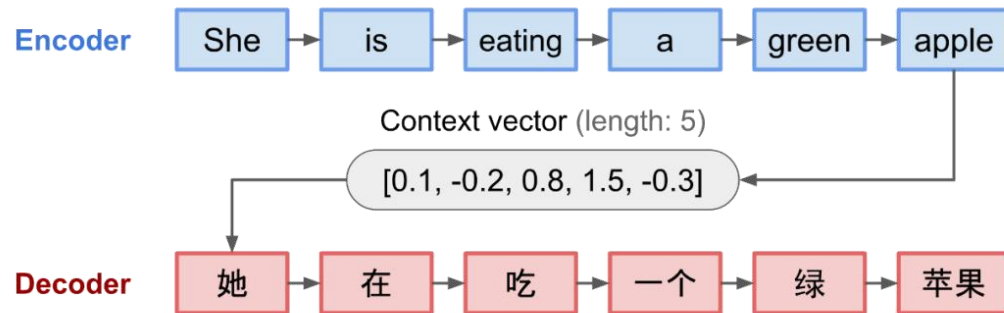
## Seq2Seq 任务示例——NMT (Neural Machine Translation)



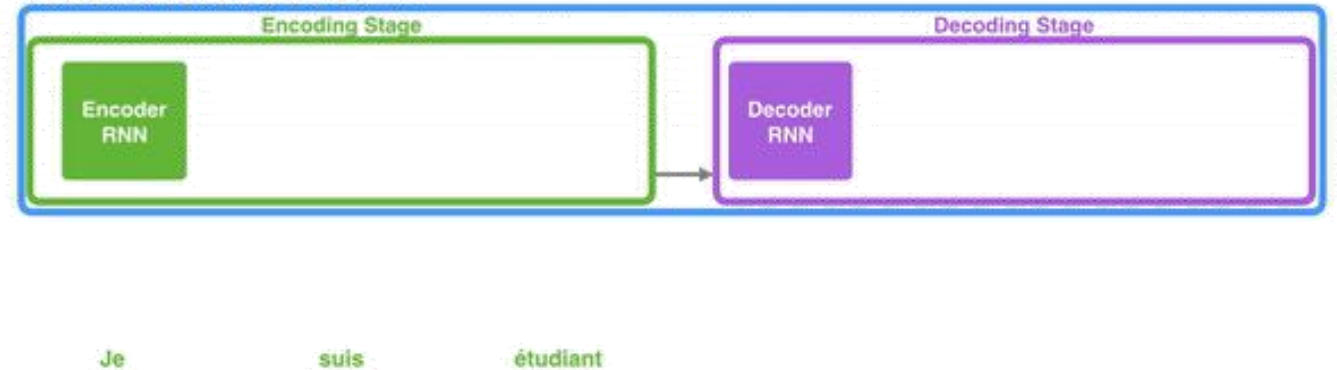## Seq2Seq模型示例——RNN (Recurrent Neural Network)



## Encoder-Decoder 框架（编码器+解码器）

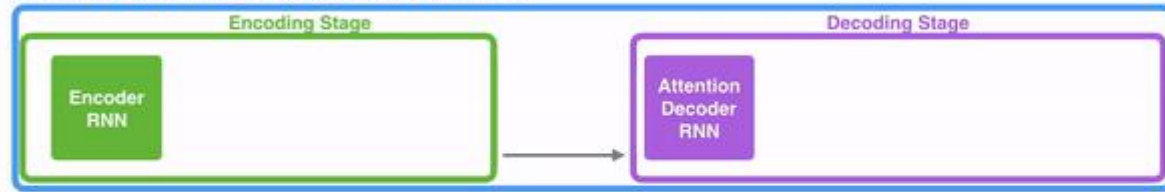# Why Transformer ?

**Seq2Seq model with Attention**



Neural Machine Translation
SEQUENCE TO SEQUENCE MODEL WITH ATTENTION

Encoding Stage    Decoding Stage

Encoder RNN    Attention Decoder RNN

Je    suis    étudiant



Neural Machine Translation
SEQUENCE TO SEQUENCE MODEL WITH ATTENTION

Encoding Stage    Attention Decoding Stage

$h_1 h_2 h_3$

$h_{init}$

<END>

4

**Attention at time step 4**

**RNN及其变体模型无法并行计算，模型效率低下！**

# Attention Is All You Need

**Ashish Vaswani***
Google Brain
avaswani@google.com

**Noam Shazeer***
Google Brain
noam@google.com

**Niki Parmar***
Google Research
nikip@google.com
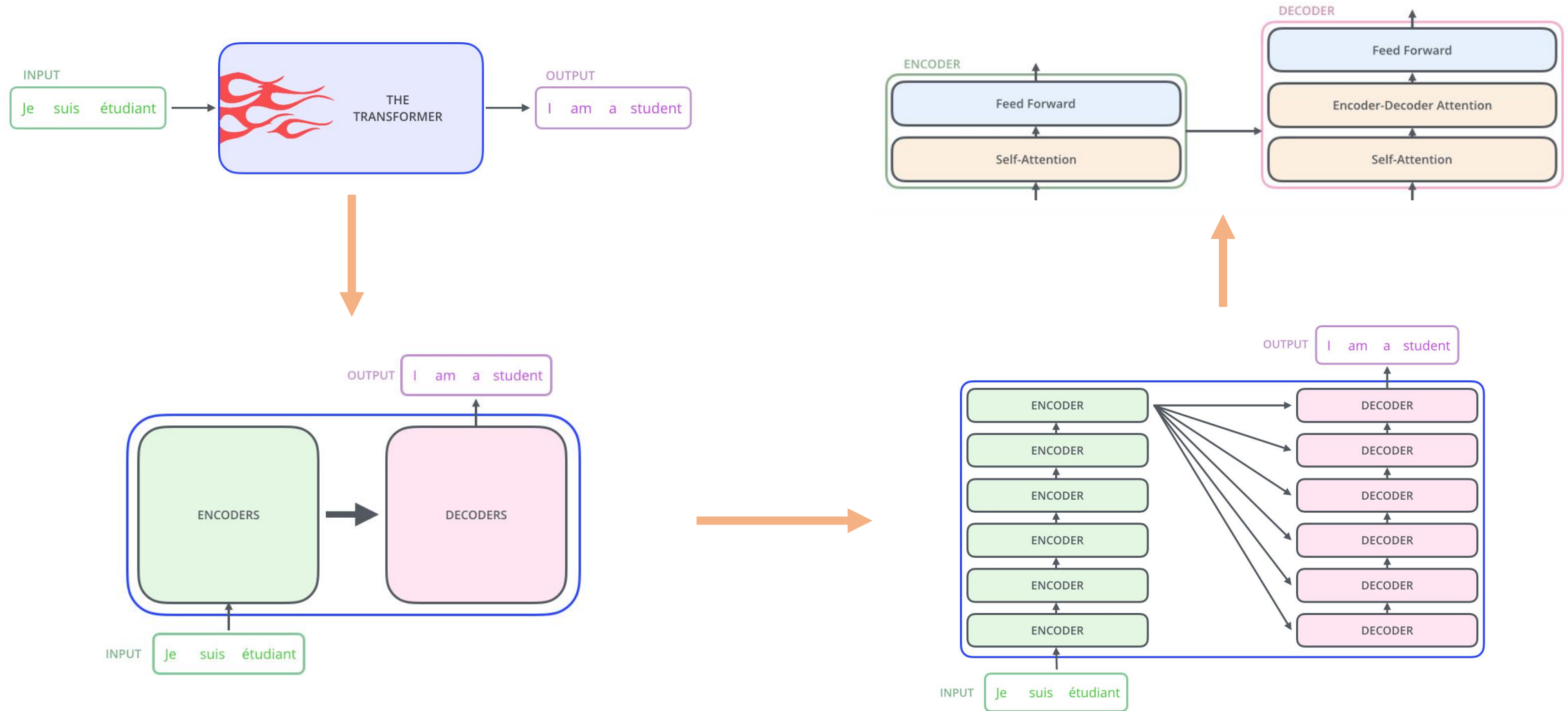
**Jakob Uszkoreit***
Google Research
usz@google.com

**Llion Jones***
Google Research
llion@google.com

**Aidan N. Gomez*** [†]
University of Toronto
aidan@cs.toronto.edu

**Łukasz Kaiser***
Google Brain
lukaszkaiser@google.com
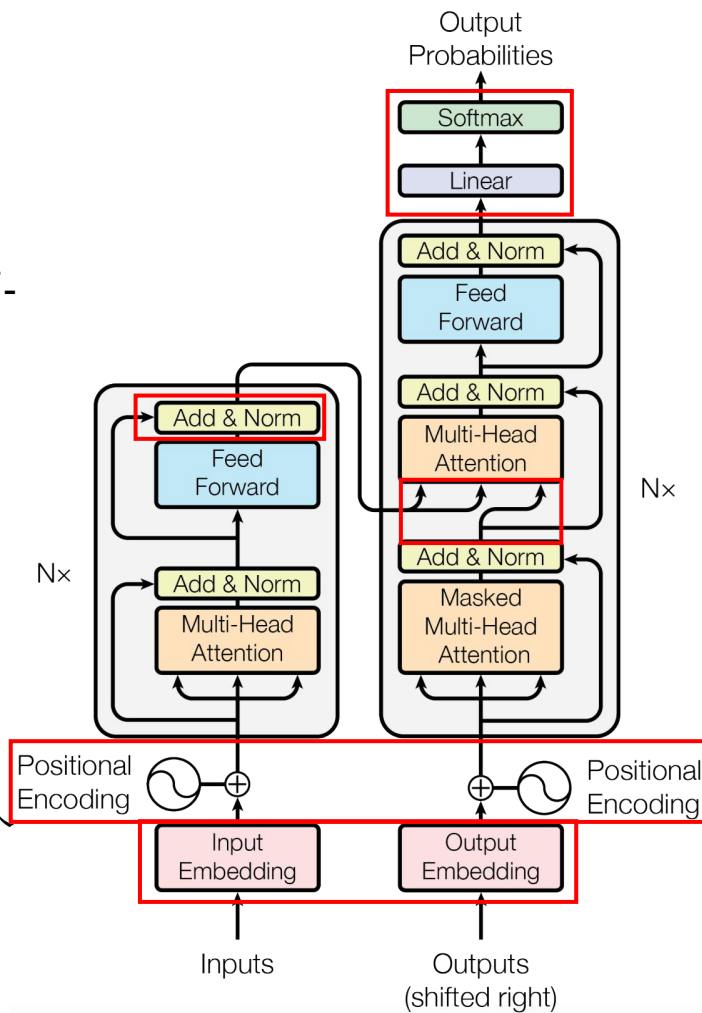
**Illia Polosukhin*** [‡]
illia.polosukhin@gmail.com

# Transformer 架构

# Transformer 架构

- 由N个block堆叠而成；
- 每个block有两层：
  - Multi-Head Attention (Self-Attention)
    + Add (Residual Connection)
    + Norm (LayerNorm)；
  - Feed Forward
    + Add (Residual Connection)
    + Norm (LayerNorm)；
- $Block_1 \sim Block_{N-1}$的输出：输入到下个Block；
- $Block_N$的输出：输入到解码器的各层中。
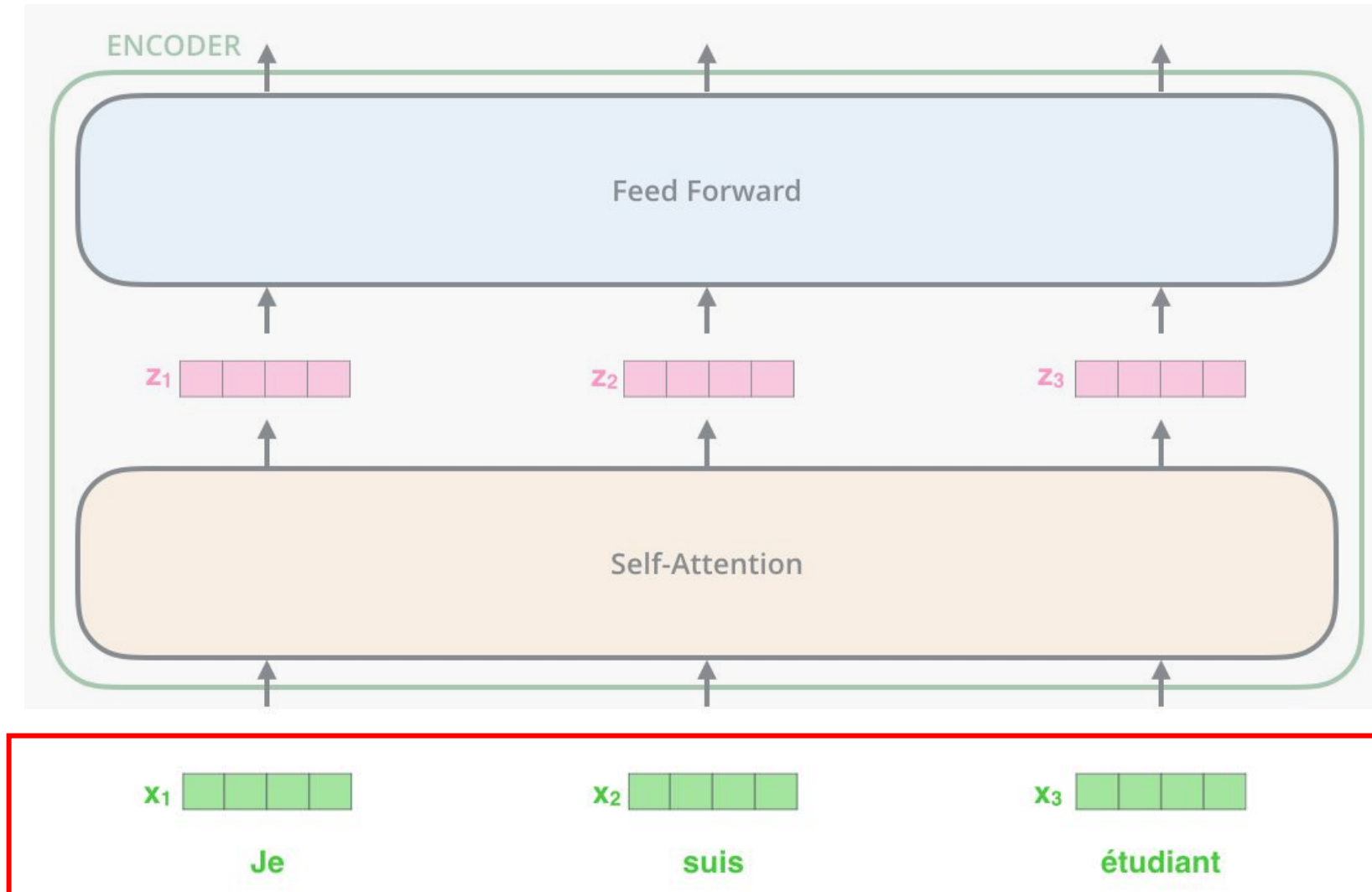


- 由N个block堆叠而成；
- 每个block有三层：
  - Masked Multi-Head Attention (Self-Attention)
    + Add (Residual Connection)
    + Norm (LayerNorm)；
  - Multi-Head Attention (Co-Attention)
    + Add (Residual Connection)
    + Norm (LayerNorm)；
  - Feed Forward
    + Add (Residual Connection)
    + Norm (LayerNorm)；
- $Block_1 \sim Block_{N-1}$的输出：输入到下个Block；
- $Block_N$的输出：输入到后续的Linear层中。

# Transformer 工作流程

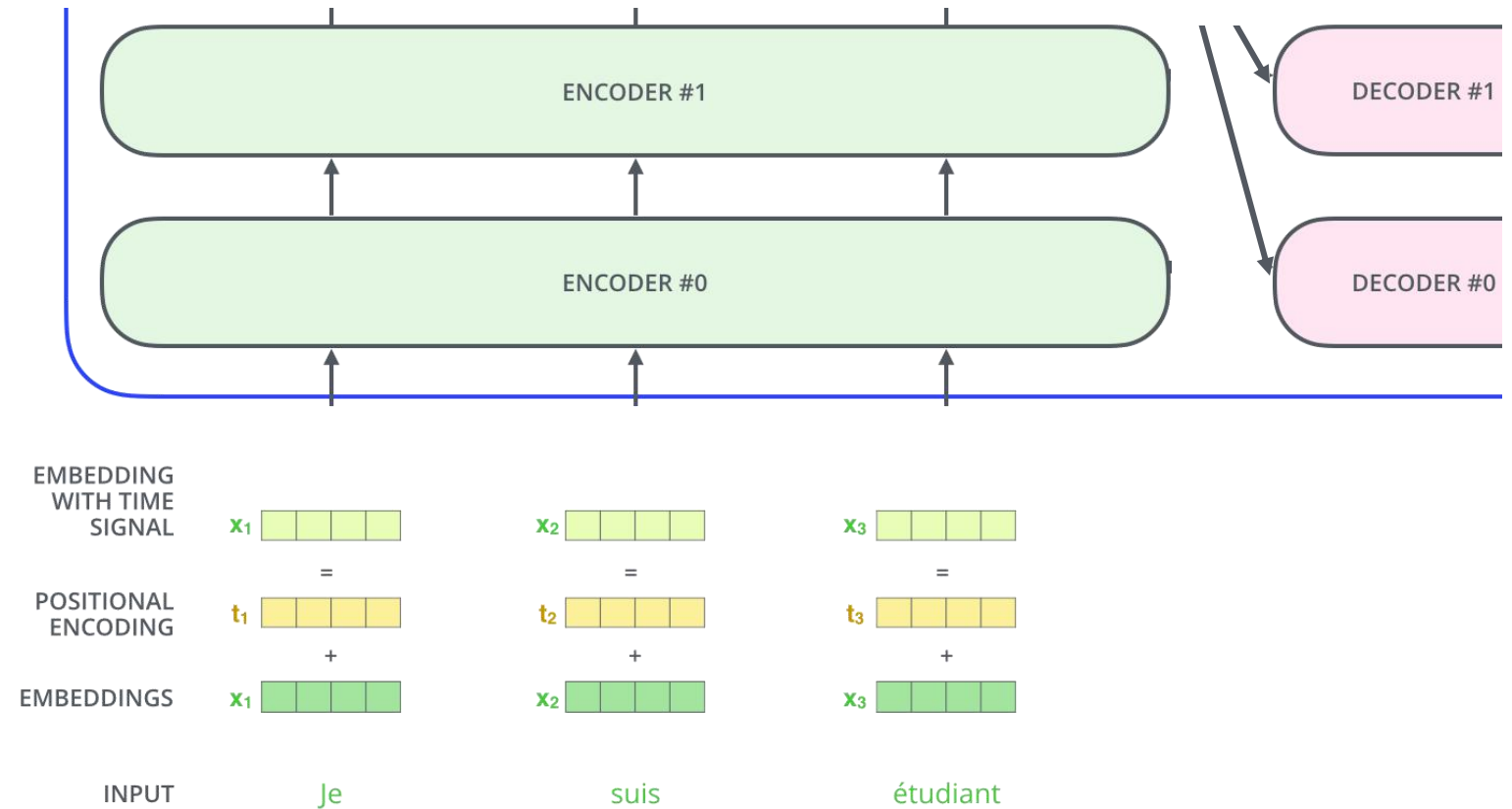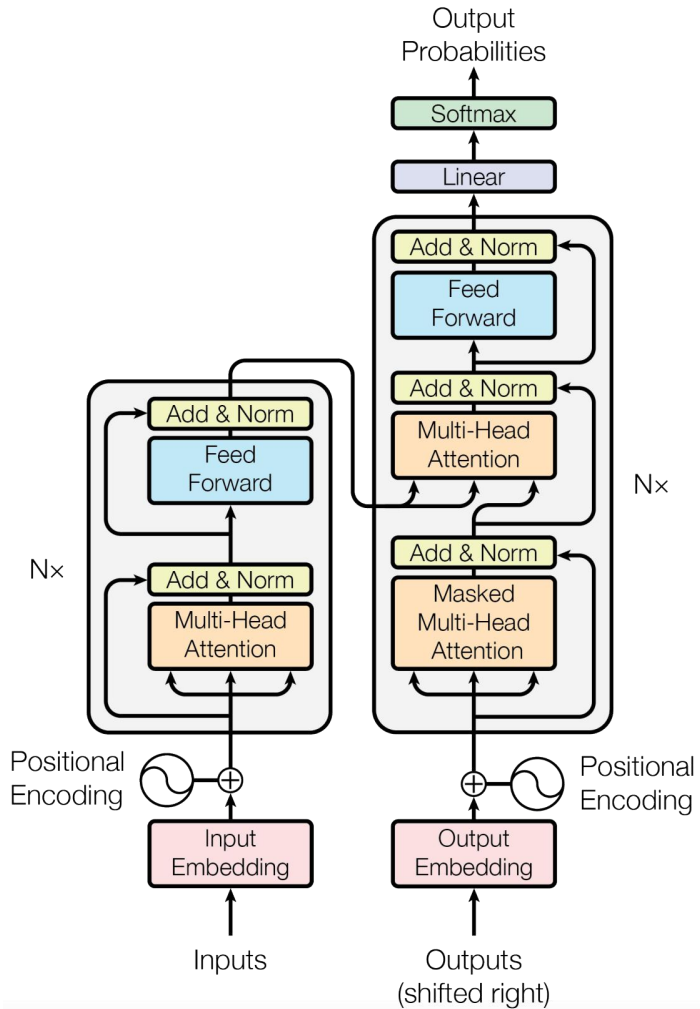**Word**

**Embedding**

## 位置编码（**Positional Encoding**）

# Transformer 工作流程

## 位置编码（**Positional Encoding**）

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

pos $\in [0, \text{max\_sequence\_length})$

$$T = \frac{2\pi}{\omega} = 2\pi * 10000^{\frac{2i}{d_{\text{model}}}}$$

I $\in [0, \frac{d_{\text{model}}}{2})$

# **Transformer** 工作流程

编码过程



ENCODER #2

ENCODER #1

$r_1$ | | | |

$r_2$ | | | |

Feed Forward Neural Network

Feed Forward Neural Network

$z_1$ | | | |

$z_2$ | | | |

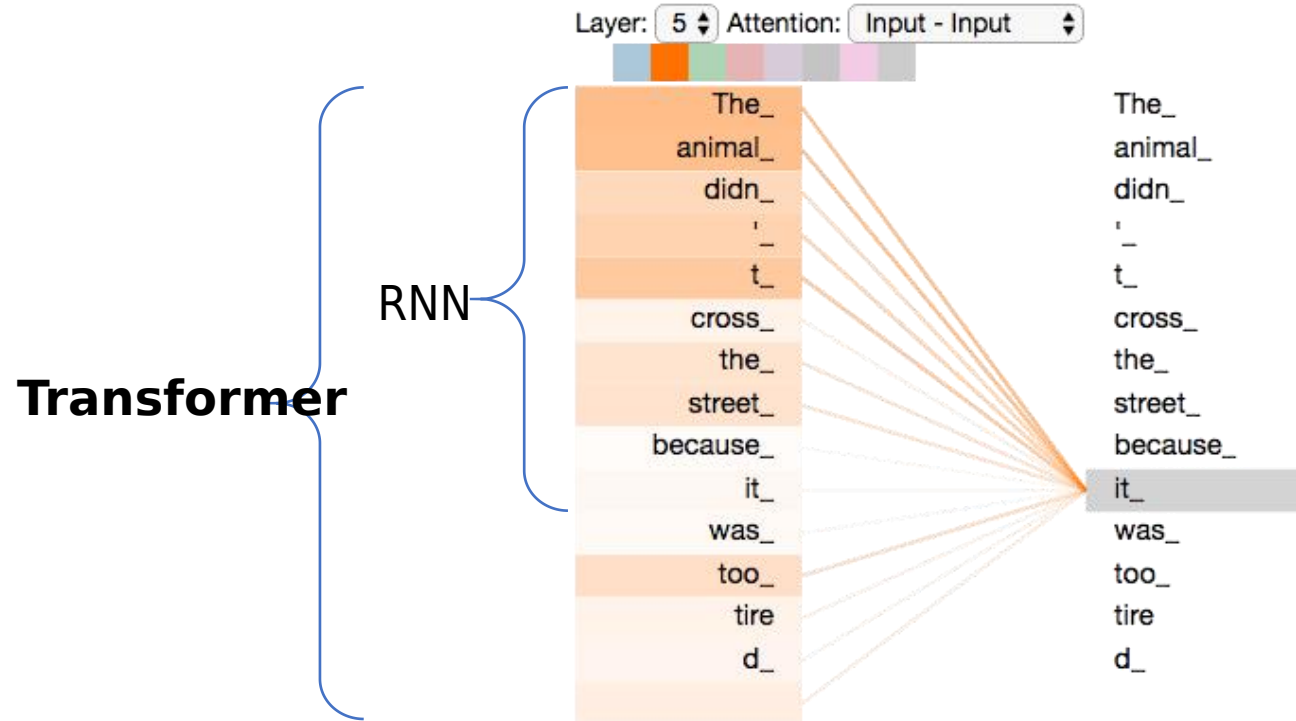Self-Attention

$x_1$ | | | |

$x_2$ | | | |

**Thinking**

**Machines**

# Transformer 工作流程

编码过程 —— **Self-Attention**（宏观）

**The animal didn't cross the street because it was too tired.**

# Transformer 工作流程

第一步：生成 **Q、K、V**，辅助计算注意力机制



$X_1$ 与 $W^Q$ 权重矩阵相乘得到 $q_1$ ,就是与这个单词相关的查询向量。通过这种方式，为输入序列的每个单词都创建一个查询向量Q、一个键向量K和一个值向量V。

# Transformer 工作流程

编码过程 —— Self-Attention（微观）

核心公式：$Attention(Q, K, V)$

$$= softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

第二步：计算当前单词的Q与候选单词的K的点积：
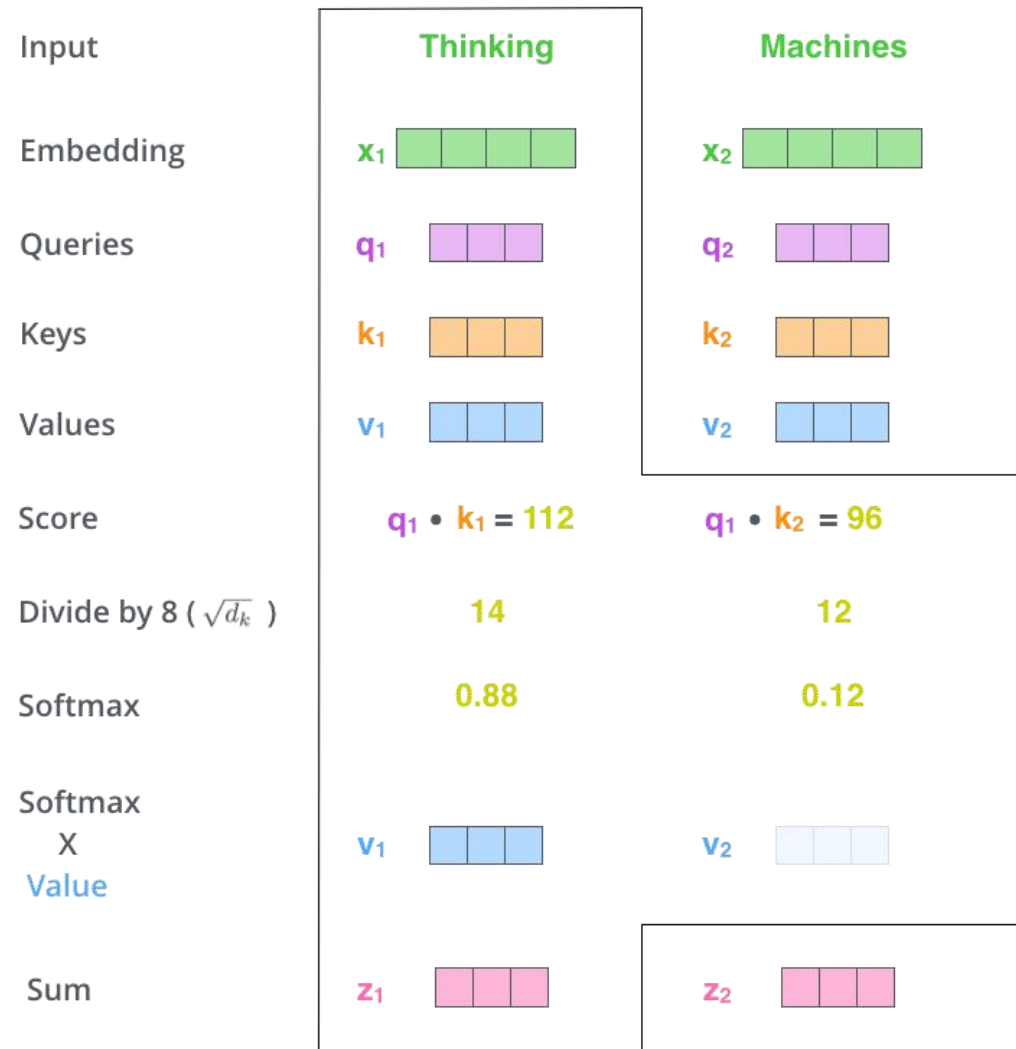
第三步：将上一步结果除以维度的平方根：

第四步：将上一步结果通过softmax函数转换：

第五步：将候选单词的每个值向量乘以softmax分数：

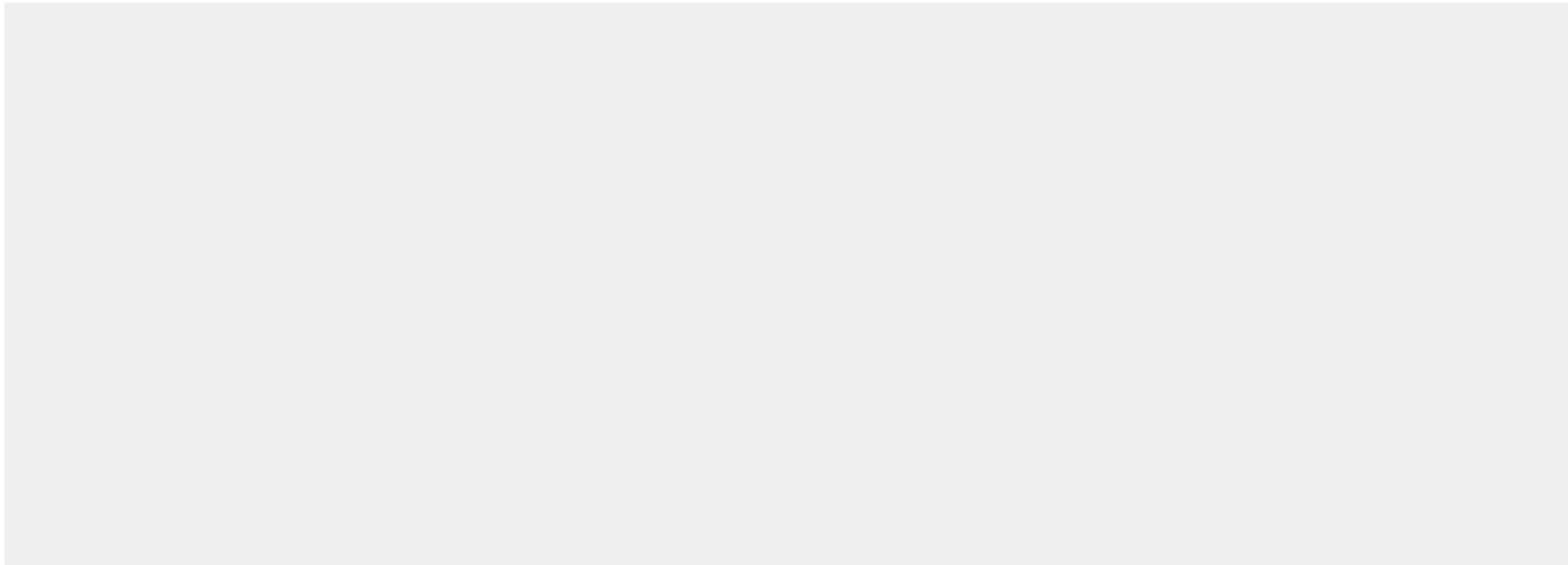第六步：对加权后的值向量求和，

即得到自注意力层在该位置的输出：

| | Thinking | Machines |
|---|---|---|
| Input | | |
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |
| Score | $q_1 \cdot k_1 = 112$ | $q_1 \cdot k_2 = 96$ |
| Divide by 8 ( $\sqrt{d_k}$ ) | 14 | 12 |
| Softmax | 0.88 | 0.12 |
| Softmax X Value | $v_1$ | $v_2$ |
| Sum | $z_1$ | $z_2$ |

# Transformer 工作流程

## 编码过程 —— Self-Attention（微观）

Self-attention



input #1

| 1 | 0 | 1 | 0 |

input #2

| 0 | 2 | 0 | 2 |

input #3

| 1 | 1 | 1 | 1 |

# Transformer 工作流程

## 编码过程 —— Self-Attention（微观）

通过矩阵运算实现自注意力机制

# Transformer 工作流程

## 编码过程 —— Multi-Head Attention

# Transformer 工作流程

## 编码过程 —— Multi-Head Attention

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

$W^O$

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

$Z$

=

ENCODER

Feed Forward

$z_1$      $z_2$      $z_3$

Self-Attention

$x_1$      $x_2$      $x_3$

Je        suis        étudiant

编码过程 —— **Multi-Head**

**Attention**

1) This is our
input sentence*

Thinking
Machines

# Transformer 工作流程

编码过程 —— **Multi-Head**

**Attention**

One-Head                     Two-Heads                     All-Heads

# Transformer 工作流程

## Padding 操作

**X**

**X**：Thinking Machines

**X** 的维度：[sequence_length, embedding_dimension]

**Actually...**

**X**：Thinking Machines _____ (seq_len: 2)
A Tale of Two Cities _____ (seq_len: 5)
Science and Art _____ (seq_len: 3)
the Art of Motorcycle Maintenance (seq_len: 5)

**X** 的维度：[batch_size, **max_sequence_length**, embedding_dimension]

max_sequence_length

**X**：

batch_size

### Scaled dot-product attention

MatMul

SoftMax

Mask (opt.)

Scale

MatMul

Q  K  V

$y = e^x$

Softmax函数: $\sigma(z_i) = \dfrac{e^{z_i}}{\sum_{j=1}^{K} e^{z_j}}$

$e^0 = 1$

$e^{-\infty} \to 0$

**Padding Mask**

| 0 | 0 | 1 | 1 | 1 |
|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 1 | 1 |
| 0 | 0 | 0 | 0 | 0 |

→

| 0 | 0 | -inf | -inf | -inf |
|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | -inf | -inf |
| 0 | 0 | 0 | 0 | 0 |

# Transformer 工作流程

## Add & Norm



**Add**

$$X_{embedding} + \text{Self Attention}(Q, K, V)$$

**Layer Norm**

$$\mu_j = \frac{1}{m} \sum_{i=1}^{m} x_{ij}$$

$$\sigma_j^2 = \frac{1}{m} \sum_{i=1}^{m} (x_{ij} - \mu_j)^2$$

$$\text{LayerNorm}(x) = \frac{x_{ij} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

**Visualization**

# Transformer 工作流程

## 编码过程

$$X_{hidden} = X_{attention} + X_{hidden}$$
$$X_{hidden} = LayerNorm(X_{hidden})$$

$$X_{hidden} = Linear(ReLU(Linear(X_{attention})))$$

$$X_{attention} = X + X_{attention}$$
$$X_{attention} = LayerNorm(X_{attention})$$

$$Q = Linear(X) = XW_Q$$
$$K = Linear(X) = XW_K$$
$$V = Linear(X) = XW_V$$
$$X_{attention} = SelfAttention(Q, K, V)$$

$$X = Embedding Lookup(X) + Positional Encoding$$

# Transformer 工作流程

Decoding time step: ①2 3 4 5 6    OUTPUT

Decoding time step: 1 ②3 4 5 6    OUTPUT    I

Linear + Softmax

ENCODER

ENCODER

DECODER

DECODER

$K_{encdec}$   $V_{encdec}$

ENCODERS

DECODERS

Linear + Softmax

EMBEDDING
WITH TIME
SIGNAL

EMBEDDINGS

INPUT    Je    suis    étudiant

EMBEDDING
WITH TIME
SIGNAL

EMBEDDINGS

INPUT    Je    suis    étudiant    PREVIOUS
OUTPUTS    I

## 解码过程 —— Masked Self-Attention



RNN模型：

# Transformer 工作流程

解码过程 —— Masked Encoder-Decoder Attention



- **Encoder Multi-Head Attention:**
  - Q, K, V all from Encoders

- **Decoder Masked Multi-Head Atte**
  - Q, K, V all from Decoders

- **Decoder Multi-Head Attention:**
  - Q from Decoders
  - K, V from Encoders

## 解码过程 —— Linear & Softmax

# Transformer 训练

| WORD | a | am | I | thanks | student | <eos> |
|---|---|---|---|---|---|---|
| INDEX | 0 | 1 | 2 | 3 | 4 | 5 |

eos：end of sentence的缩写形式

输入： "je suis étudiant"，期望输出： "i am a student"

**Target Model Outputs**

Output Vocabulary:    a    am    I    thanks  student  <eos>

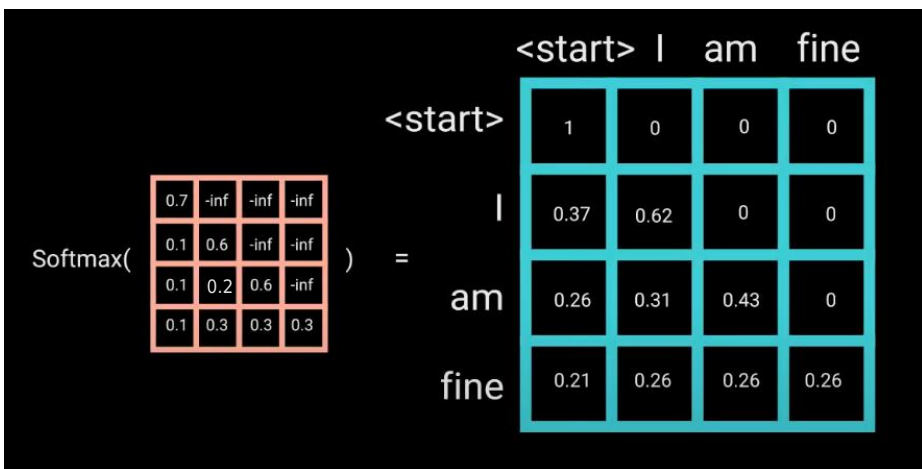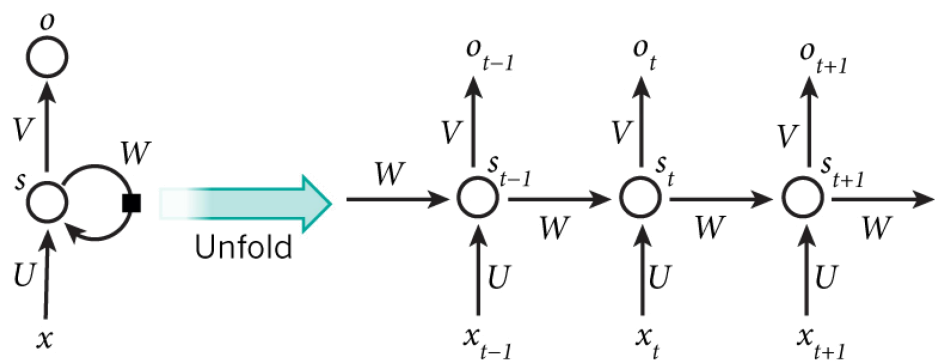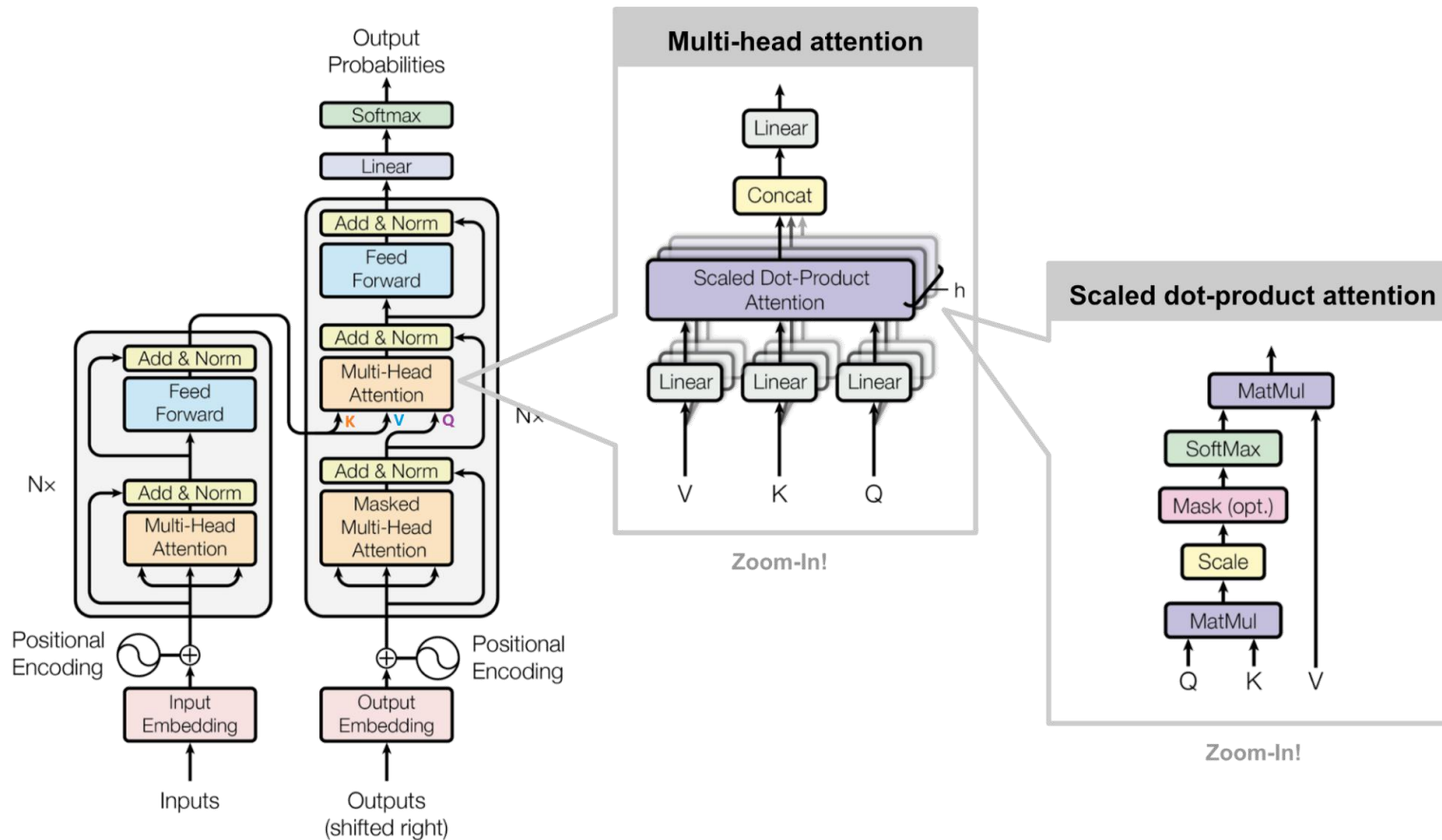| | a | am | I | thanks | student | <eos> |
|---|---|---|---|---|---|---|
| position #1 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 |
| position #2 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| position #3 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| position #4 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 |
| position #5 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 |

a    am    I    thanks  student  <eos>

**Trained Model Outputs**

Output Vocabulary:    a    am    I    thanks  student  <eos>

| | a | am | I | thanks | student | <eos> |
|---|---|---|---|---|---|---|
| position #1 | 0.01 | 0.02 | 0.93 | 0.01 | 0.03 | 0.01 |
| position #2 | 0.01 | 0.8 | 0.1 | 0.05 | 0.01 | 0.03 |
| position #3 | 0.99 | 0.001 | 0.001 | 0.001 | 0.002 | 0.001 |
| position #4 | 0.001 | 0.002 | 0.001 | 0.02 | 0.94 | 0.01 |
| position #5 | 0.01 | 0.01 | 0.001 | 0.001 | 0.001 | 0.98 |

a    am    I    thanks  student  <eos>

27

# Transformer 预测

**解码/采样策略**



| | | | |
|---|---|---|---|
| 0.01 | 0.8 | 0.1 | 0.05 | 0.01 | 0.03 |

**概率降序排列**

| 0.8 | 0.1 | 0.05 | 0.03 | 0.01 | 0.01 |
|---|---|---|---|---|---|

**缩小采样范围**

Top-k (k=3):  | 0.8 | 0.1 | 0.05 |

Top-p (p=0.95): | 0.8 | 0.1 | 0.05 | 0.03 |

**随机采样**

Top-k (k=3):  | 0.8 |

Top-p (p=0.95): | 0.1 |

Argmax Decoding（贪婪采样）
　　Greedy Search（贪心搜索）
　　Beam Search（集束搜索）
Stochastic Decoding（随机采样）
　　Temperature-controlled Stochastic Sampling
　　Top-k Sampling
　　Top-p Sampling（Nucleus Sampling）
混合采样

CTRL： $p_i = \frac{\exp(x_i/(T \cdot I(i \in g)))}{\sum_j \exp(x_j/(T \cdot I(j \in g)))}$   $I(c) = \theta$ if $c$ is True else 1

Timestep 1 2 3 4

| | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| A | 0.5 | 0.1 | 0.2 | 0.0 |
| B | 0.2 | 0.4 | 0.2 | 0.2 |
| C | 0.2 | 0.3 | 0.4 | 0.2 |
| <eos> | 0.1 | 0.2 | 0.2 | 0.6 |

$P(x|x_{1:t-1}) = \frac{exp(u_t/t)}{\sum_{t'} exp(u_{t'}/t)}$ ，其中 $t \in [0, 1)$

# 作业

实现 **Transformer** 架构，通过单步调试理解数据流动过程及维度变化。

参考代码：https://github.com/wmathor/nlp-tutorial

## 5. Model based on Transformer

- 5-1. The Transformer - **Translate**
  - Paper - Attention Is All You Need(2017)
  - Colab - Transformer_Torch.ipynb
  - bilibili - https://www.bilibili.com/video/BV1mk4y1q7eK

| Model | Example |
|---|---|
| NNLM | Predict Next Word |
| Word2Vec(Softmax) | Embedding Words and Show Graph |
| TextCNN | Sentence Classification |
| TextRNN | Predict Next Step |
| TextLSTM | Autocomplete |
| Bi-LSTM | Predict Next Word in Long Sentence |
| Seq2Seq | Change Word |
| Seq2Seq with Attention | Translate |
| Bi-LSTM with Attention | Binary Sentiment Classification |
| Transformer | Translate |
| Greedy Decoder Transformer | Translate |
| BERT | how to train |

# 参考资料

- Attention is All you Need (acm.org)

- The Illustrated Transformer – Jay Alammar – Visualizing machine learning one concept at a time. (jalammar.github.io)

- The Annotated Transformer (harvard.edu)

- wmathor.com

- Transformer详解_数学家是我理想的博客-CSDN博客

- 图解Transformer（完整版）_龙心尘的博客-CSDN博客

- 如何优雅地编码文本中的位置信息？三种positioanl encoding方法简述 (qq.com)

- NLP基础模型和注意力机制_开始King的博客-CSDN博客

- 文本生成自回归解码策略总结_自回归解码器_Meilinger_的博客-CSDN博客

# Q & A