

Transformer模型

《大语言模型》编写团队：李军毅

大模型神经网络的奠基之作



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

Attention is all you need

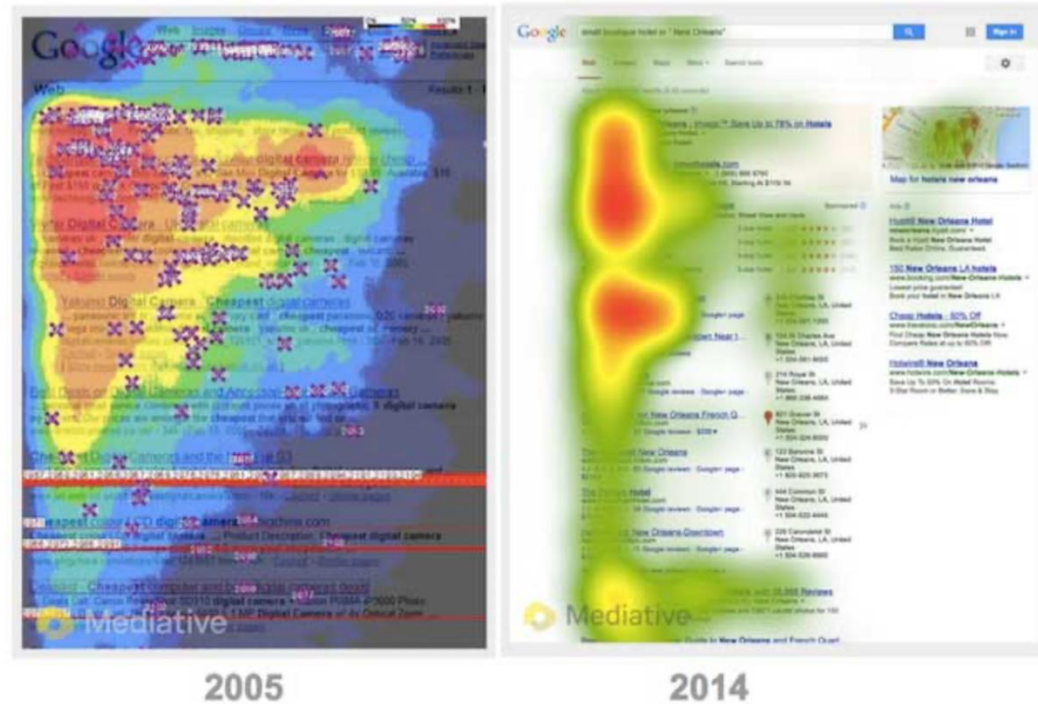
[\[PDF\] neurips.cc](#)

[A Vaswani](#), [N Shazeer](#), [N Parmar](#)... - *Advances in neural ...*, 2017 - [proceedings.neurips.cc](#)

... **Jakob** proposed replacing RNNs with self-**attention** and started the effort to evaluate this idea. **Ashish**, with **Illia**, designed and ... **Noam** proposed scaled dot-product **attention**, multi-head ...

☆ 保存 ↻ 引用 被引用次数: 153082 相关文章 所有 91 个版本 ↗

Why Attention?



Source: The Evolution of Google Search Results Pages, Mediative, 2014

在**搜索场景**中，人们的目光往往会更加关注**左上角**的三角区域
(即第一条搜索结果的位置)

A [similar study](#) conducted by Mediative in 2005 found that users tended to focus their gaze on the "Golden Triangle" — the top-left corner of a SERP where the first result was usually displayed.

The 2014 results were strikingly different from those a decade ago. The Golden Triangle has all but disappeared; instead, users tend scan more vertically and to vary their focus depending on what they are searching for.

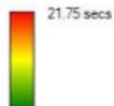
Why Attention?



在阅读过程中

Why Attention?

Media: Diapers-01.jpg
Time: 00:00:00.000 - 00:00:06.033
Participant filter: All



Extra gentle for the most sensitive skin.

So gentle on ultra sensitive skin, add the chemicals and moisture of a diaper and you have diaper rash.

Baby Wipes's unique high-absorbency natural-blend cotton lining provides cotton-soft, extra thick, gel-free protection for your baby's sensitive skin. The chlorine-free materials and absorbent polymers is non-toxic and non-irritating. Clinically tested and pediatrician recommended for babies with allergies and sensitive skin.

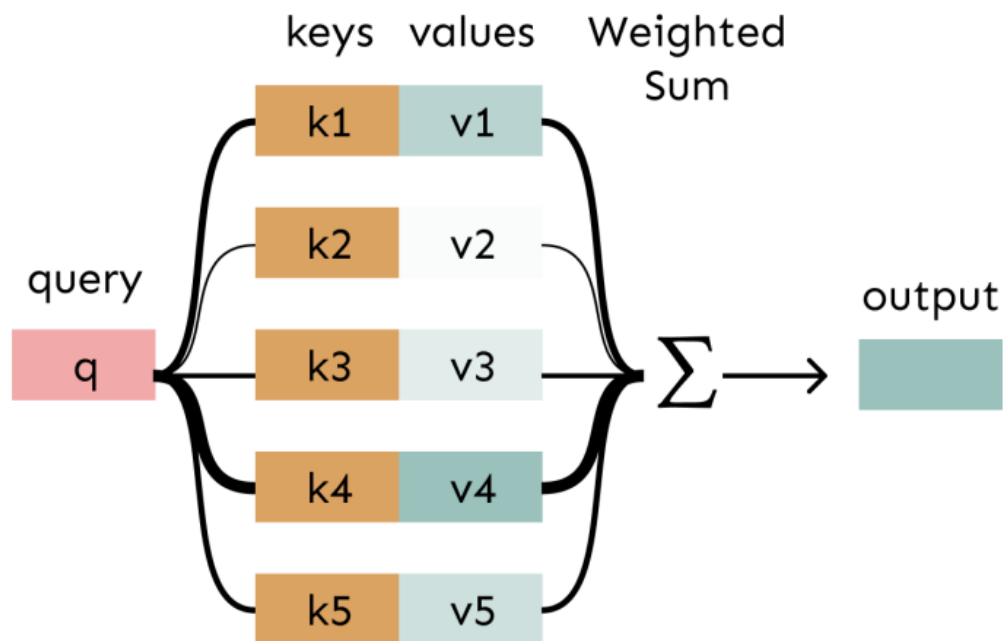


Wipes by™

If you are not satisfied with the baby leakage protection, you will get your money back. Read more about our leakfree guarantee at www.baby.com

在浏览图文信息时

➤ 可以视为一种基于相似度的查表



➤ 第一步：计算query与key相似度

$$e_{ij} = q_i^T k_j$$

➤ 第二步：相似度规范化

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{j'} \exp(e_{ij'})}$$

➤ 第三步：对value加权求和

$$o_i = \sum_j \alpha_{ij} v_i$$

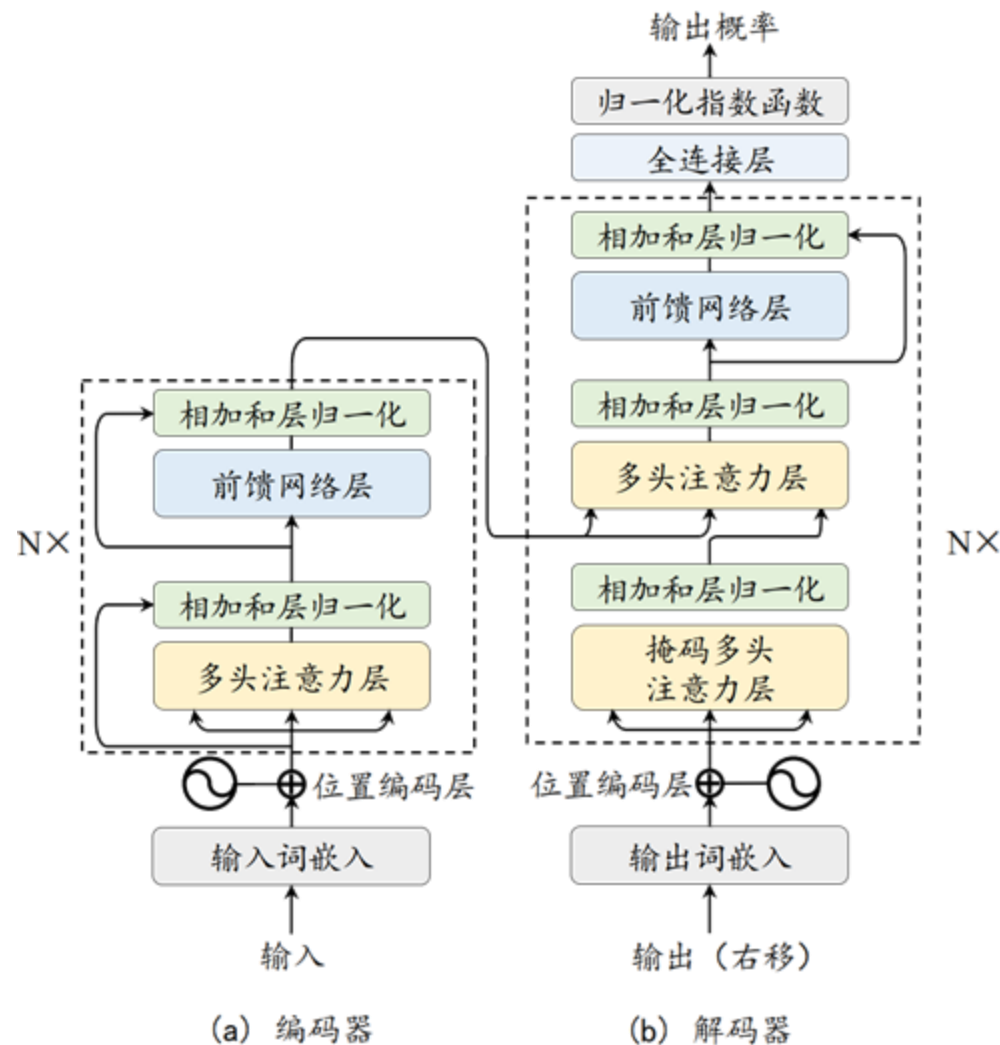
Transformer模型

➤ 核心模块：注意力

- Transformer 完全抛弃传统的CNN 和 RNN，
整个网络结构完全由注意力机制组成

➤ 编码器-解码器结构

- 编码器将输入序列变换为隐藏层特征
- 解码器将隐藏层特征变换为输出序列



Transformer模型

➤ 编码器：将输入变换为隐藏层特征

➤ N 个堆叠的编码器层

➤ 多头注意力

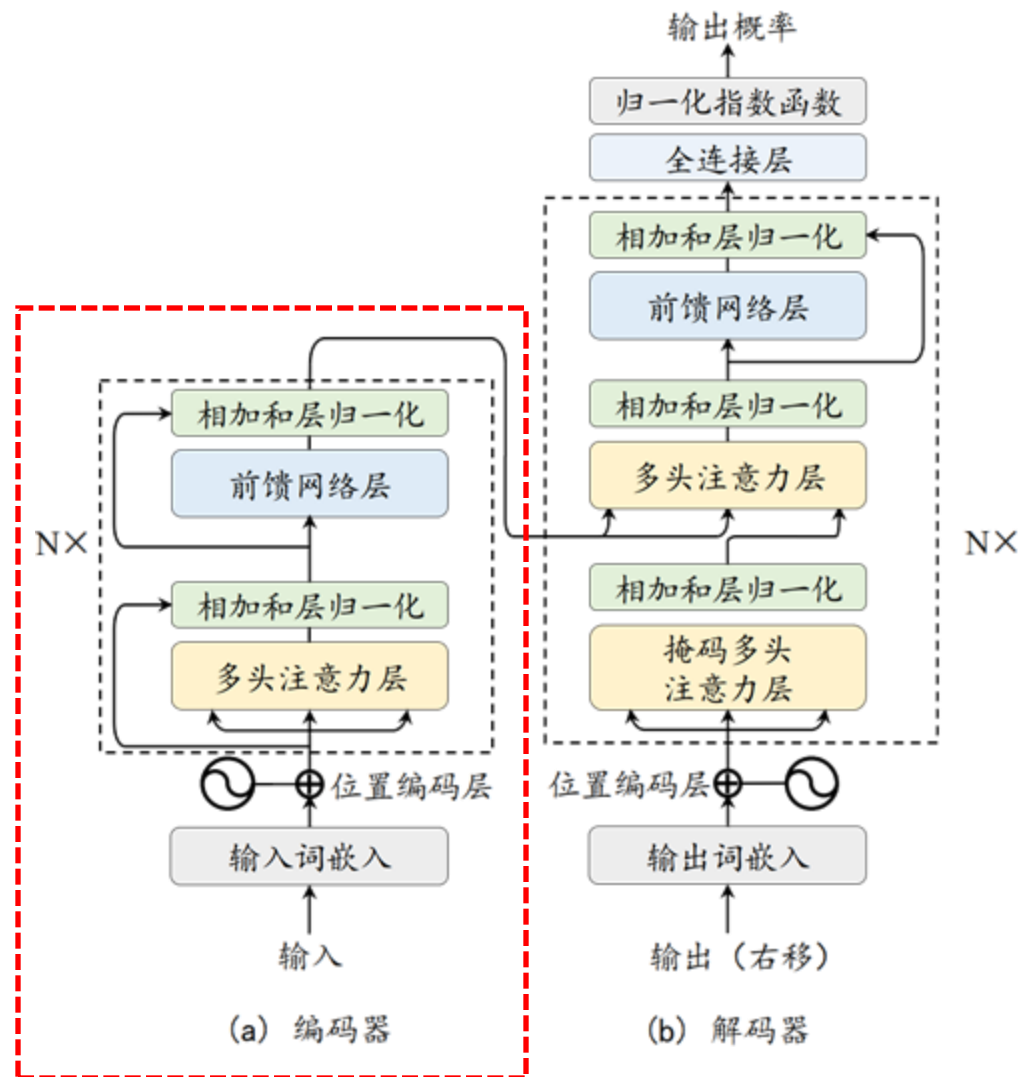
➤ 前馈网络

➤ 残差连接和层归一化

$$X'_l = \text{LayerNorm}(\text{MHA}(X_{l-1}) + X_{l-1}),$$

$$X_l = \text{LayerNorm}(\text{FFN}(X'_l) + X'_l),$$

X_{l-1} ：编码器第 $l-1$ 层的输出



Transformer模型

➤ 解码器：将隐藏层特征变换为自然语言序列

➤ N 个堆叠的解码器层

➤ (掩码) 多头注意力

➤ 前馈网络

➤ 残差连接和层归一化

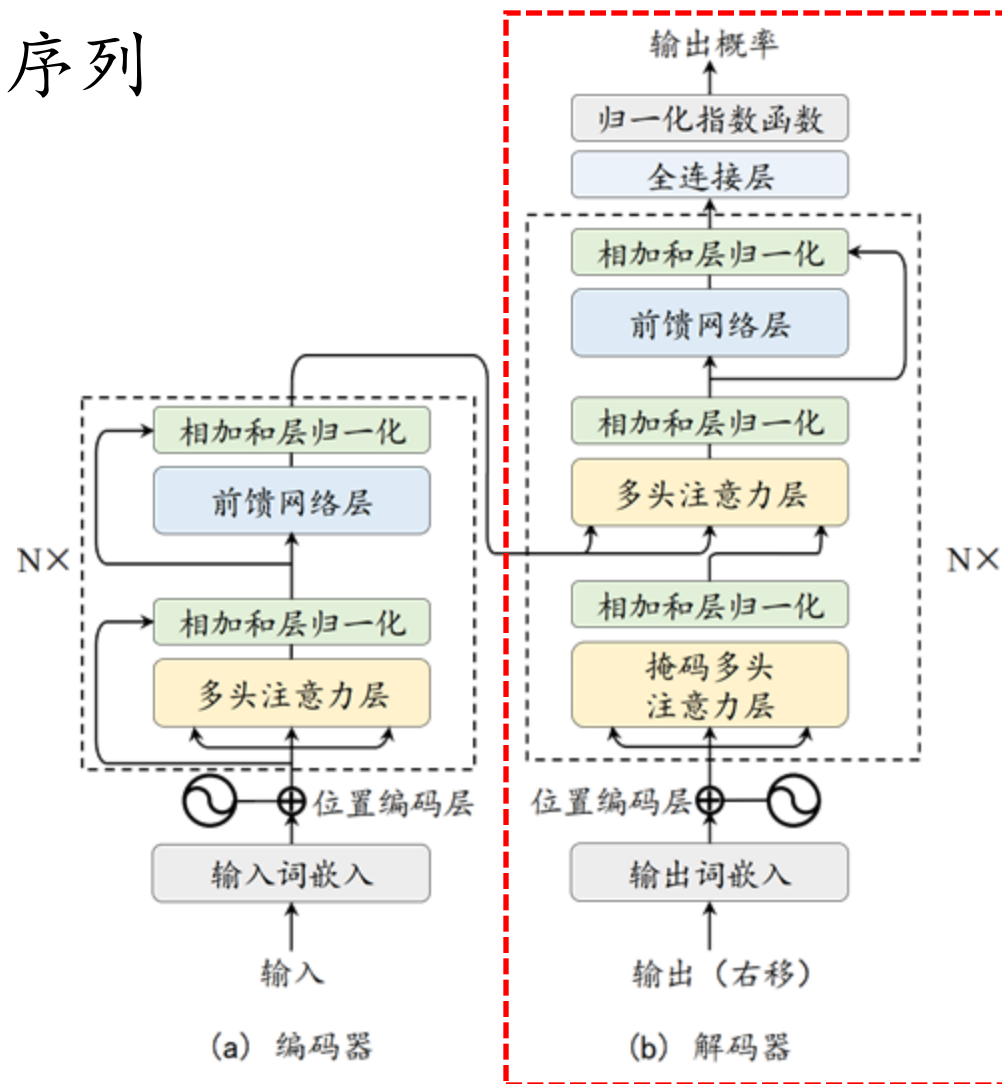
$$Y'_l = \text{LayerNorm}(\text{MaskedMHA}(Y_{l-1}) + Y_{l-1}),$$

$$Y''_l = \text{LayerNorm}(\text{CrossMHA}(Y'_l, X_L) + Y'_l),$$

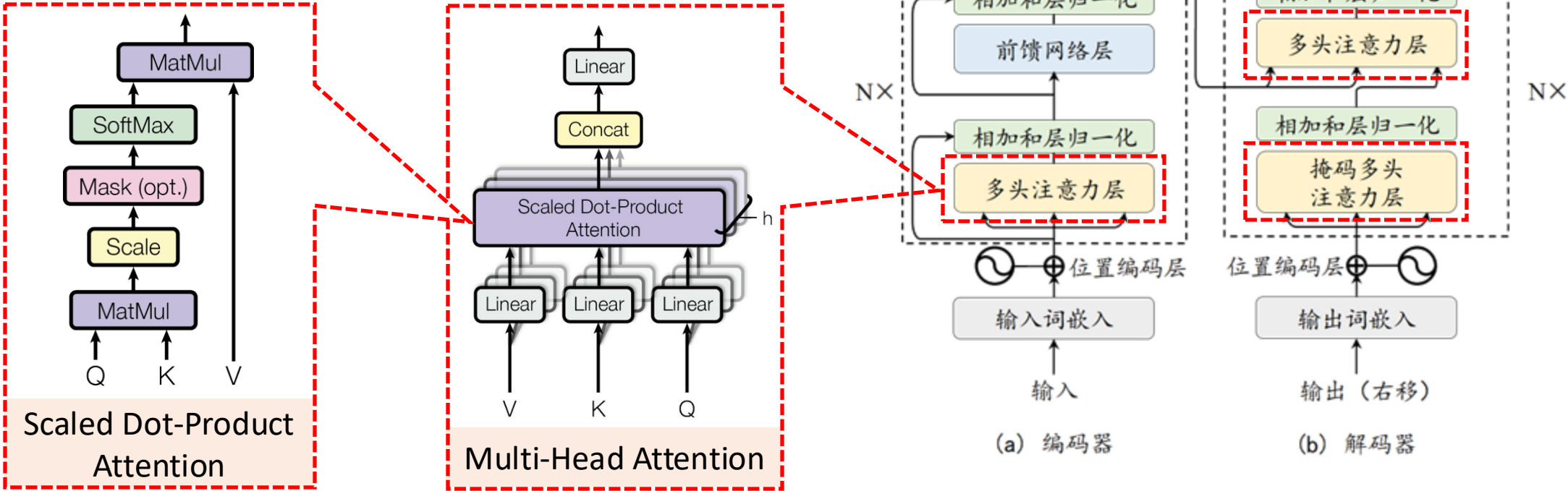
$$Y_l = \text{LayerNorm}(\text{FFN}(Y''_l) + Y''_l),$$

Y_{l-1} : 解码器第 $l-1$ 层的输出

X_L : 编码器第 L 层的输出



- 多头注意力层：建模任意距离词元间的语义关系
- Scaled Dot-Product Attention (单注意力头计算)
- Multi-Head Attention (多注意力头拼接)



➤ 多头注意力层

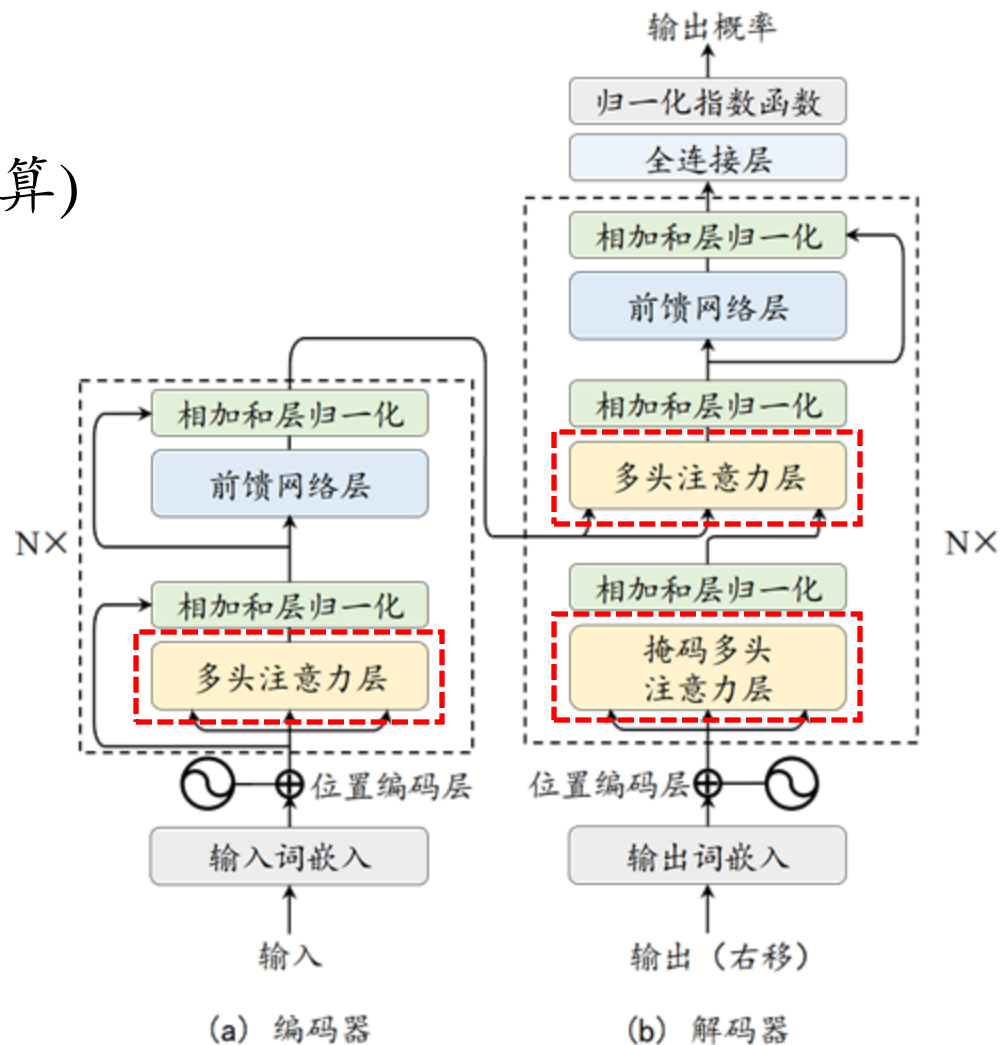
➤ Scaled Dot-Product Attention (单注意力头计算)

➤ Multi-Head Attention (多注意力头拼接)

思考：为什么拆分再合并？

➤ 将注意力分为多个头形成多个子空间，可以让模型去关注不同方面信息，最后再将各个方面的信息综合起来

➤ 多次注意力计算综合的结果类比 CNN 中同时使用多个卷积核的作用



➤ 多头注意力层

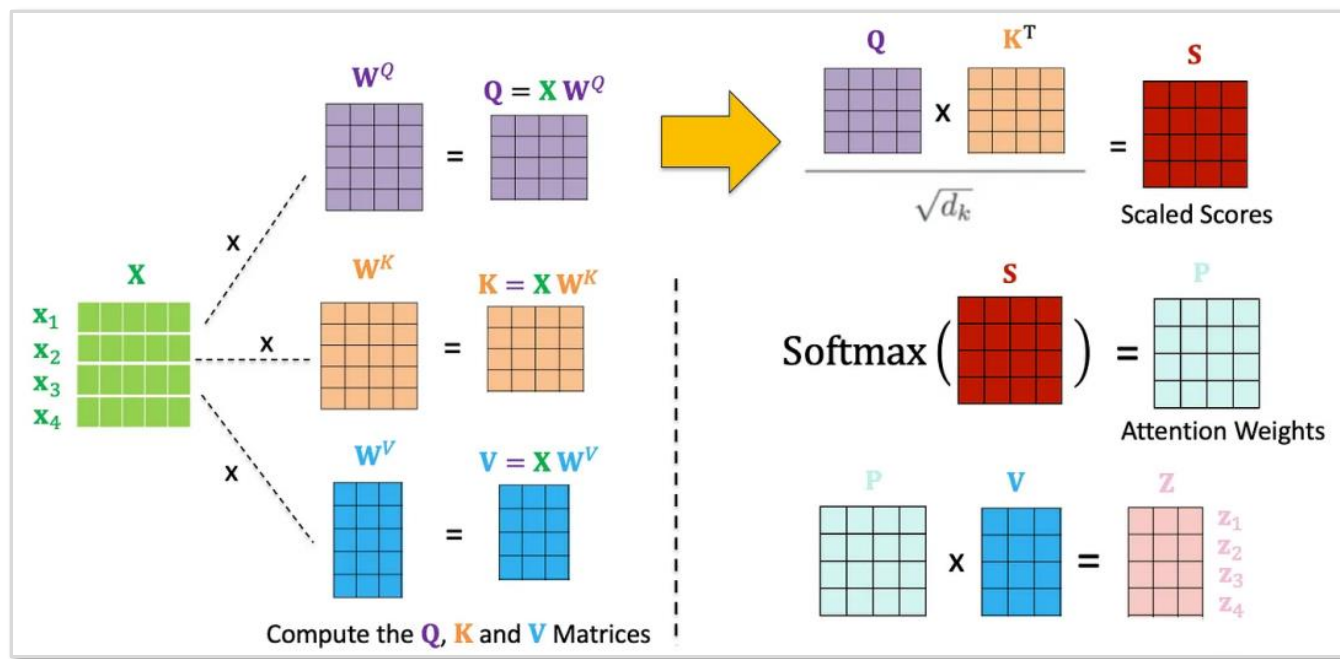
➤ Scaled Dot-Product Attention

➤ 将输入 X 映射为 Q 、 K 、 V 矩阵

$$Q = XW^Q, \quad K = XW^K, \quad V = XW^V,$$

➤ 注意力计算

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



➤ 多头注意力层

➤ Scaled Dot-Product Attention

➤ 将输入 X 映射为 Q 、 K 、 V 矩阵

$$Q = XW^Q, \quad K = XW^K, \quad V = XW^V,$$

➤ 注意力计算

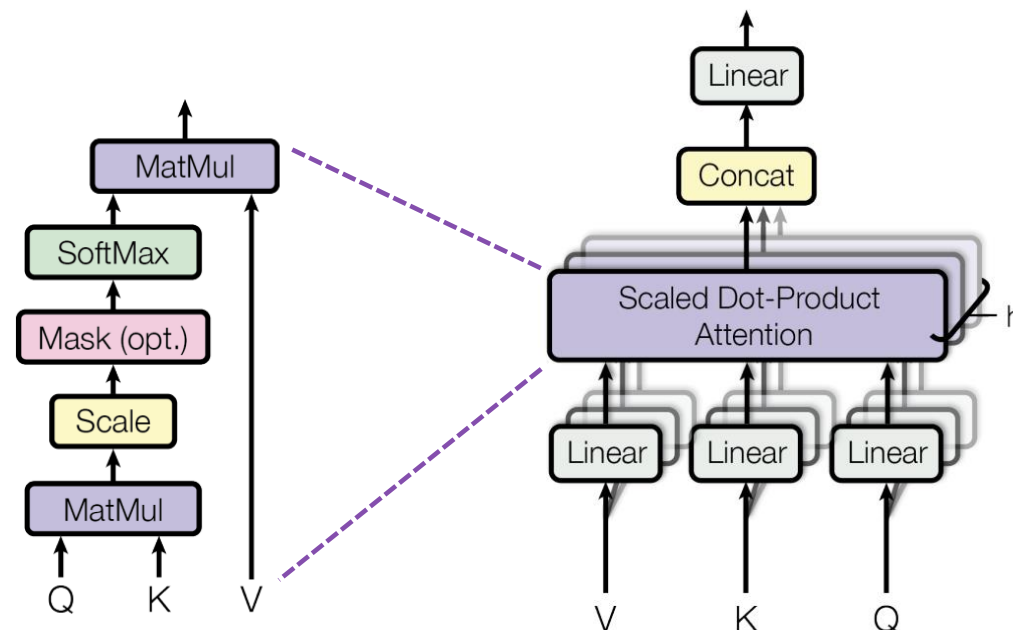
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{D}}\right)V$$

➤ Multi-Head Attention

➤ 拼接多个注意力头

$$\text{MHA} = \text{Concat}(\text{head}_1, \dots, \text{head}_N)W^O,$$

$$\text{head}_n = \text{Attention}(XW_n^Q, XW_n^K, XW_n^V).$$



Scaled Dot-Product Attention

Multi-Head Attention

➤ 思考

➤ Q, K, V 在编码器、解码器中各是什么？

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

Transformer模型

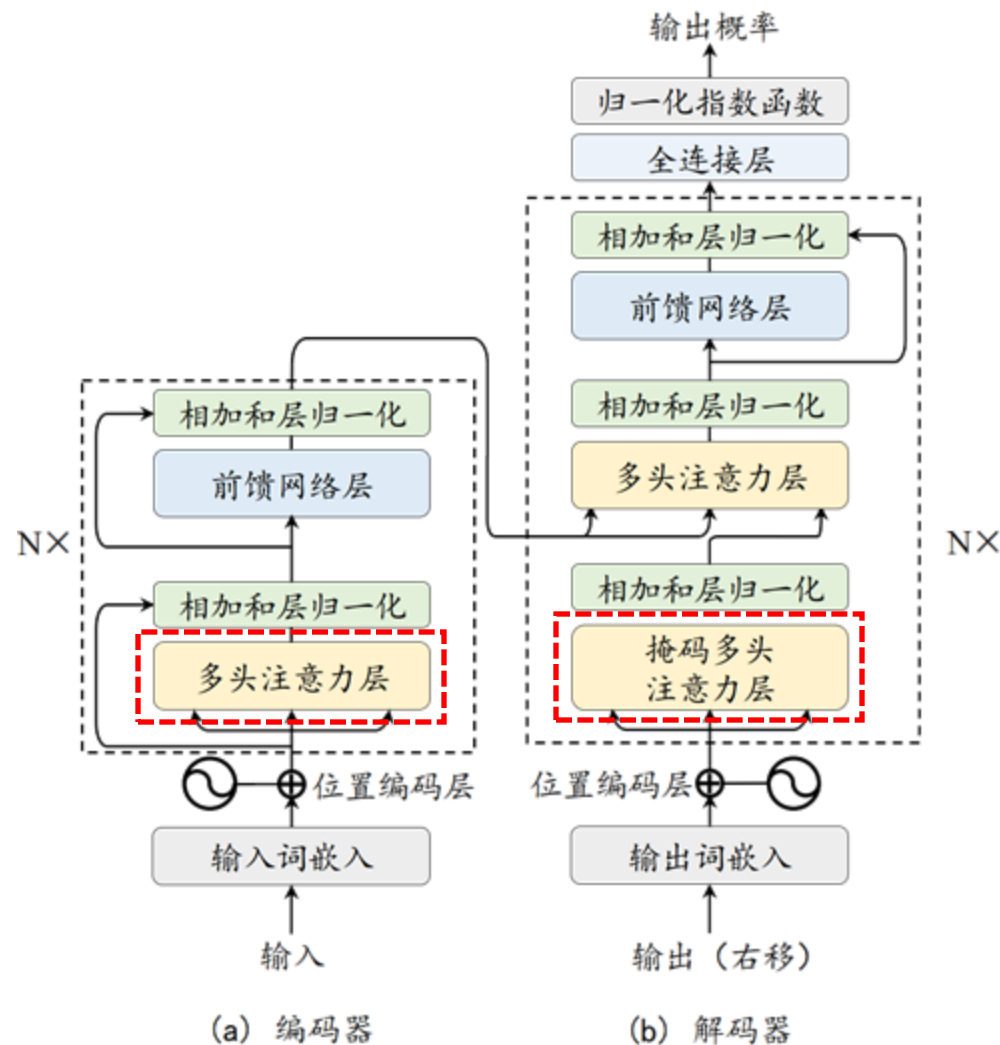
➤ 思考

➤ Q, K, V 在编码器、解码器中各是什么？

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

➤ 在编码器、解码器中

Q, K, V 相同，均为自身前一层输出
(名称Self-Attention的由来)



Transformer模型

➤ 思考

➤ Q, K, V 在编码器、解码器中各是什么？

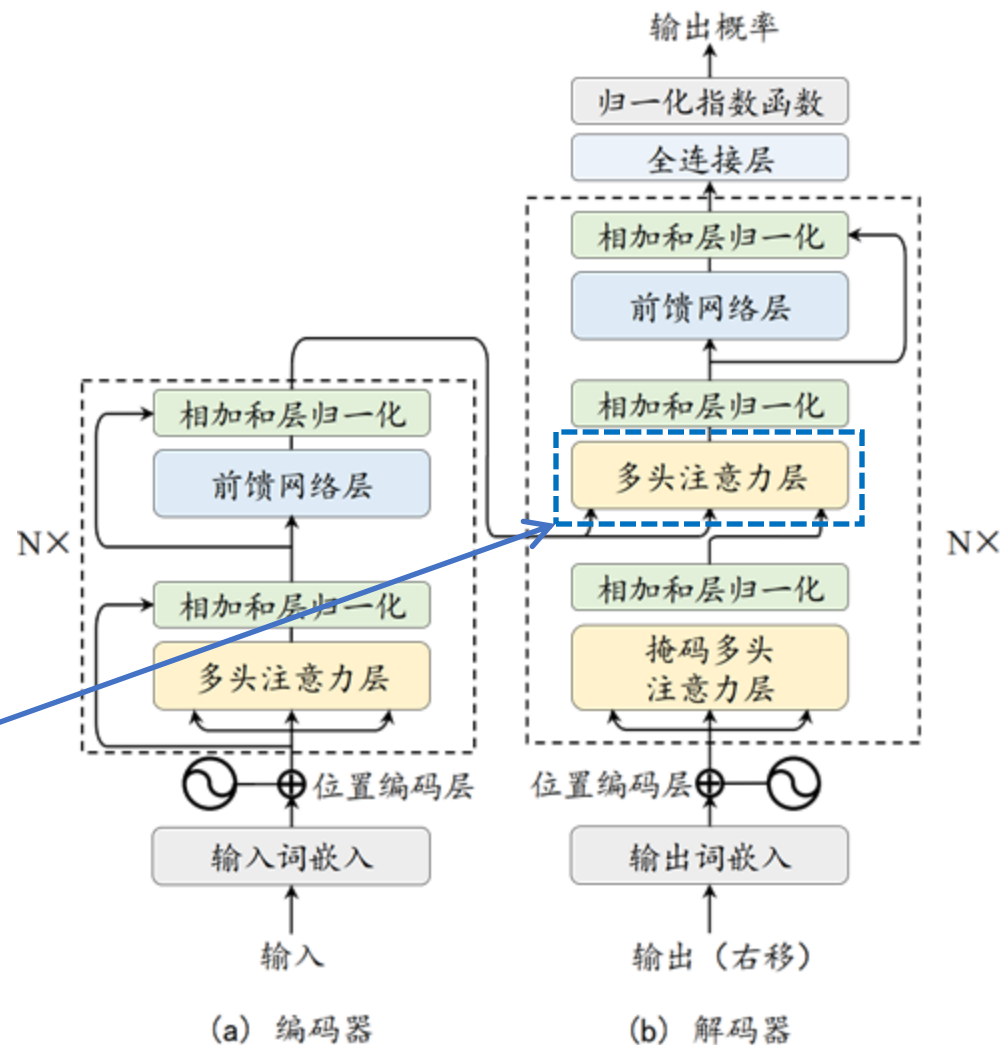
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{D}}\right)V$$

➤ 在编码器、解码器中

Q, K, V 相同，均为自身前一层输出

➤ 唯一不同

Q 来自前一层输出， K, V 是编码器输出



➤ 思考

➤ 为什么要用点积? 有其他候选吗?

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

➤ 思考

➤ 为什么要用点积？有其他候选吗？

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

The two most commonly used attention functions are **additive attention** [2], and **dot-product** (multiplicative) attention. Dot-product attention is identical to our algorithm, except for the scaling factor of $\frac{1}{\sqrt{d_k}}$. Additive attention computes the compatibility function using a feed-forward network with a single hidden layer. While the two are similar in theoretical complexity, **dot-product attention is much faster and more space-efficient in practice, since it can be implemented using highly optimized matrix multiplication code.**

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})}$$
$$e_{ij} = v_a^\top \tanh(W_a s_{i-1} + U_a h_j),$$

➤ 思考

➤ 为什么要scale, 即除以 \sqrt{D} (D 为隐藏层特征维度)?

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

➤ 思考

➤ 为什么要scale, 即除以 \sqrt{D} (D 为隐藏层特征维度)?

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

d_k 也是 D While for small values of d_k the two mechanisms perform similarly, additive attention outperforms dot product attention without scaling for larger values of d_k [3]. We suspect that for large values of d_k , the dot products grow large in magnitude, pushing the softmax function into regions where it has extremely small gradients⁴. To counteract this effect, we scale the dot products by $\frac{1}{\sqrt{d_k}}$.

⁴To illustrate why the dot products get large, assume that the components of q and k are independent random variables with mean 0 and variance 1. Then their dot product, $q \cdot k = \sum_{i=1}^{d_k} q_i k_i$, has mean 0 and variance d_k .

Transformer模型

➤ 输入编码

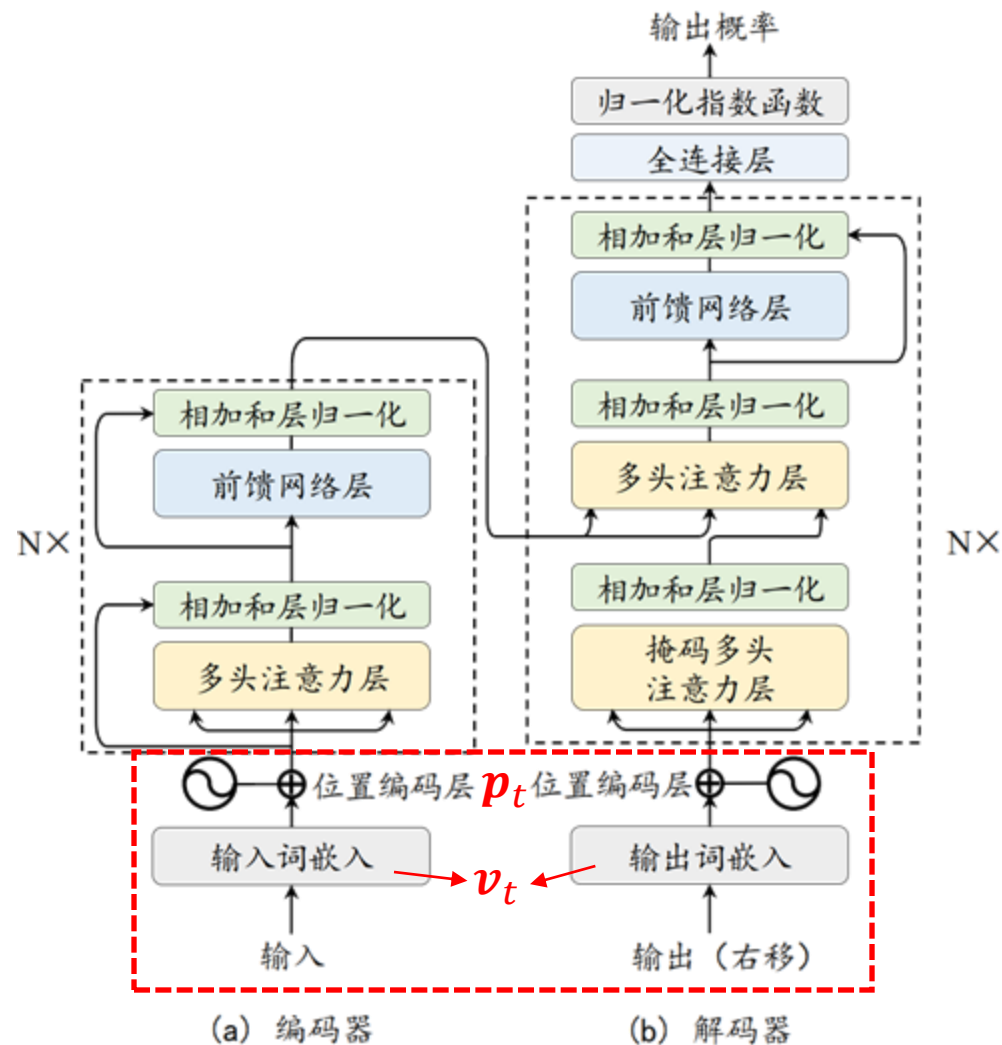
$$\mathbf{x}_t = \mathbf{v}_t + \mathbf{p}_t$$

➤ \mathbf{v}_t : 词嵌入 (语义信息)

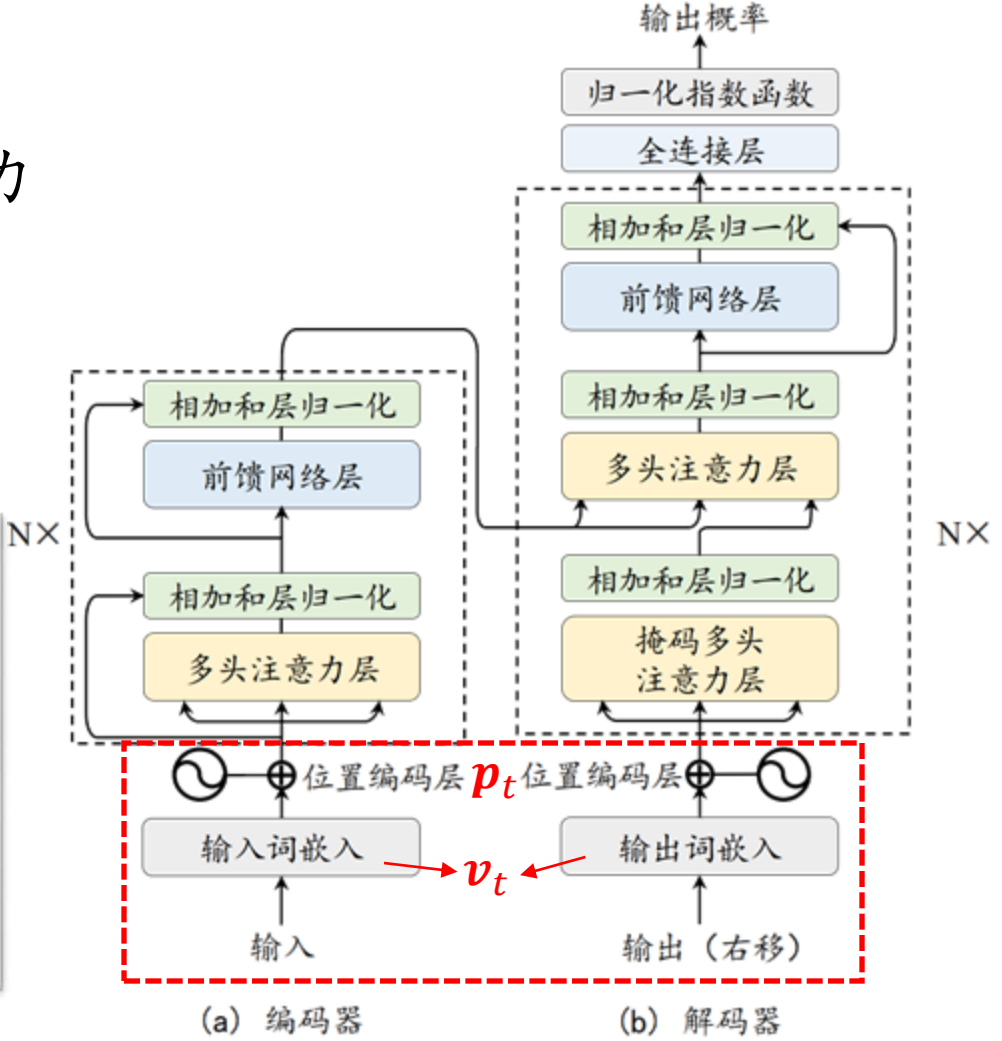
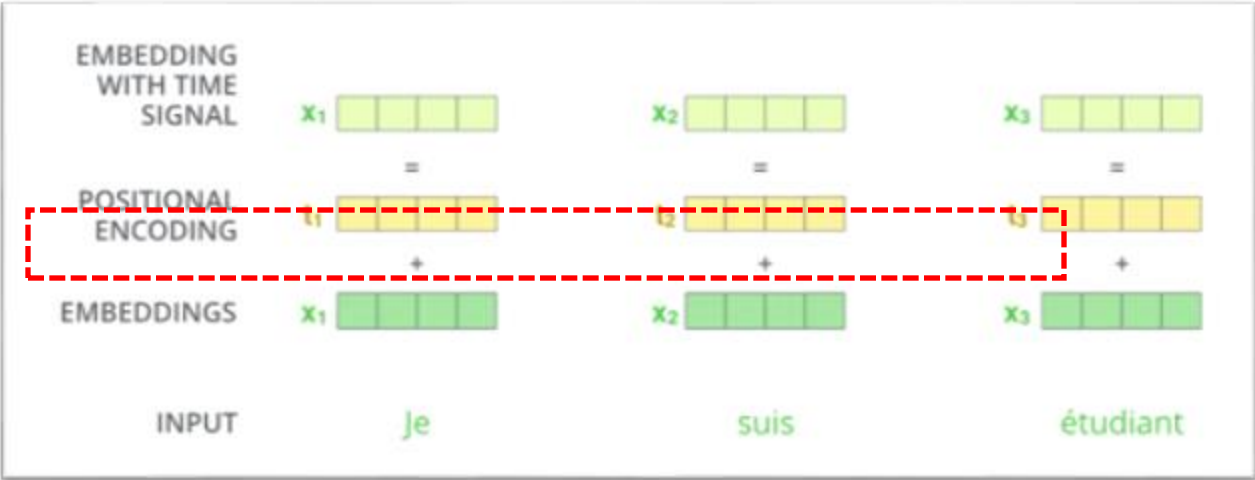
➤ 词元通过嵌入模块映射为词向量

➤ \mathbf{p}_t : 位置编码 (位置信息)

➤ 根据词元的绝对位置分配位置向量



- 位置编码
 - Transformer 本身不具备建模顺序序列的能力
 - 在输入词嵌入中引入位置编码，
将绝对或相对位置信息注入



➤ 位置编码

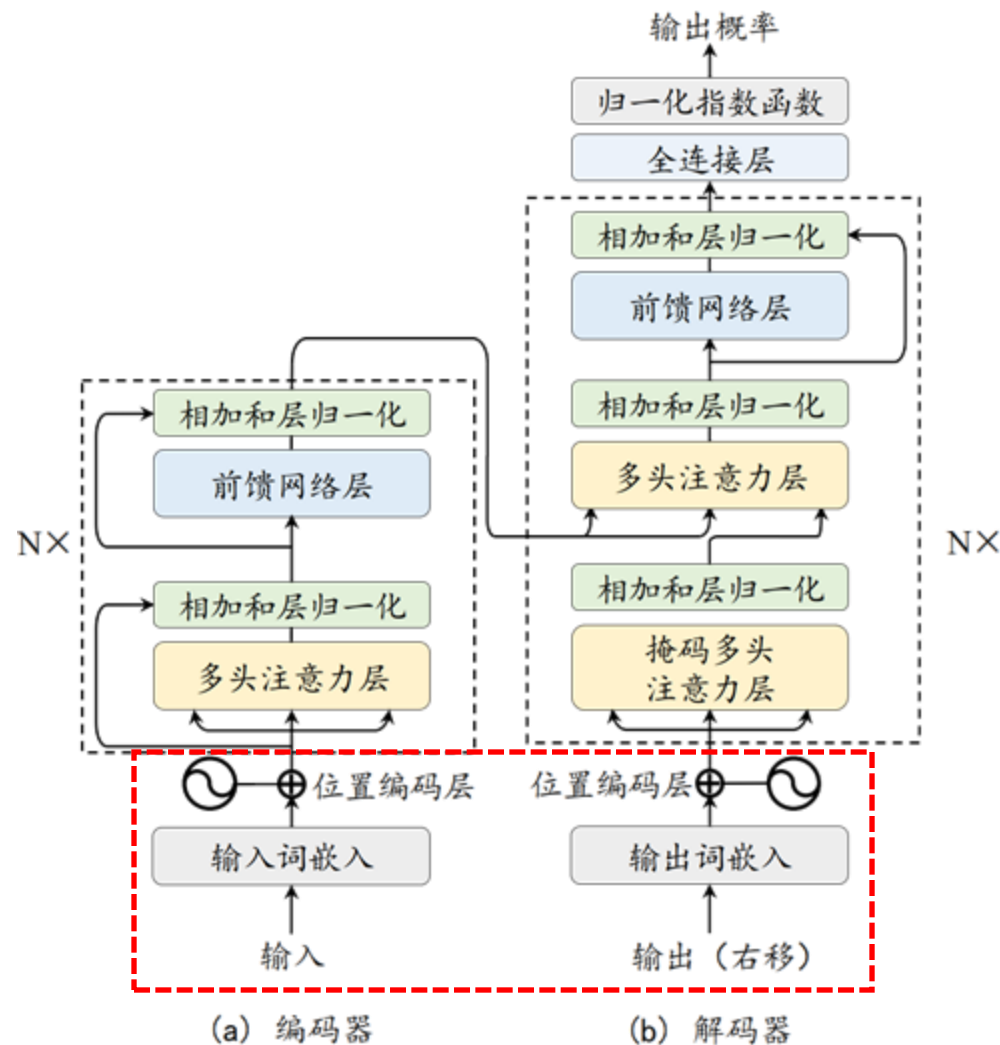
➤ 可以通过训练学习得到，也可以事先指定

➤ 论文使用 \sin 和 \cos 函数线性变换得到

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

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

pos 表示一句话中单词的位置， i 是词嵌入维度序号， d_{model} 是词嵌入维度



➤ 位置编码

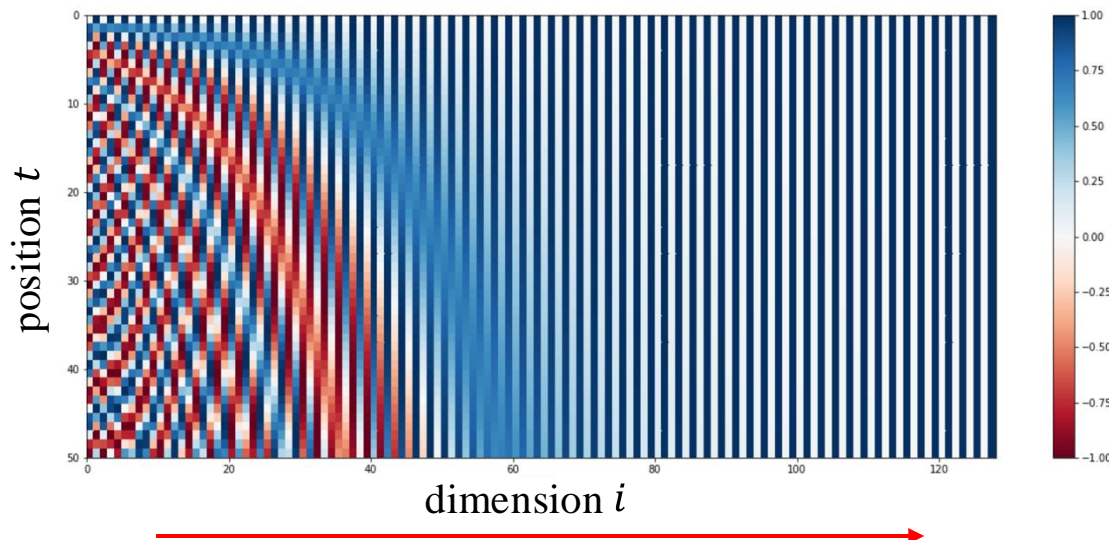
➤ 将 sin 和 cos 函数转换得到

$$p_t^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k + 1 \end{cases} \quad \omega_k = \frac{1}{10000^{2k/d}}$$

词嵌入维度序号 i 变大 (k 变大),
 ω_k 越小, 频率 $\frac{\omega_k}{2\pi}$ 也越来越小

➤ 一个词的位置编码表示如下

$$p_t = \begin{bmatrix} \sin(w_1 t) \\ \cos(w_2 t) \\ \dots \\ \sin(w_{d/2} t) \\ \cos(w_{d/2} t) \end{bmatrix}_{d \times 1}$$



随着词嵌入维度序号 i 增大, 该
序号下的数值变化频率
指数级下降

➤ 前馈网络层：学习复杂的函数关系和特征

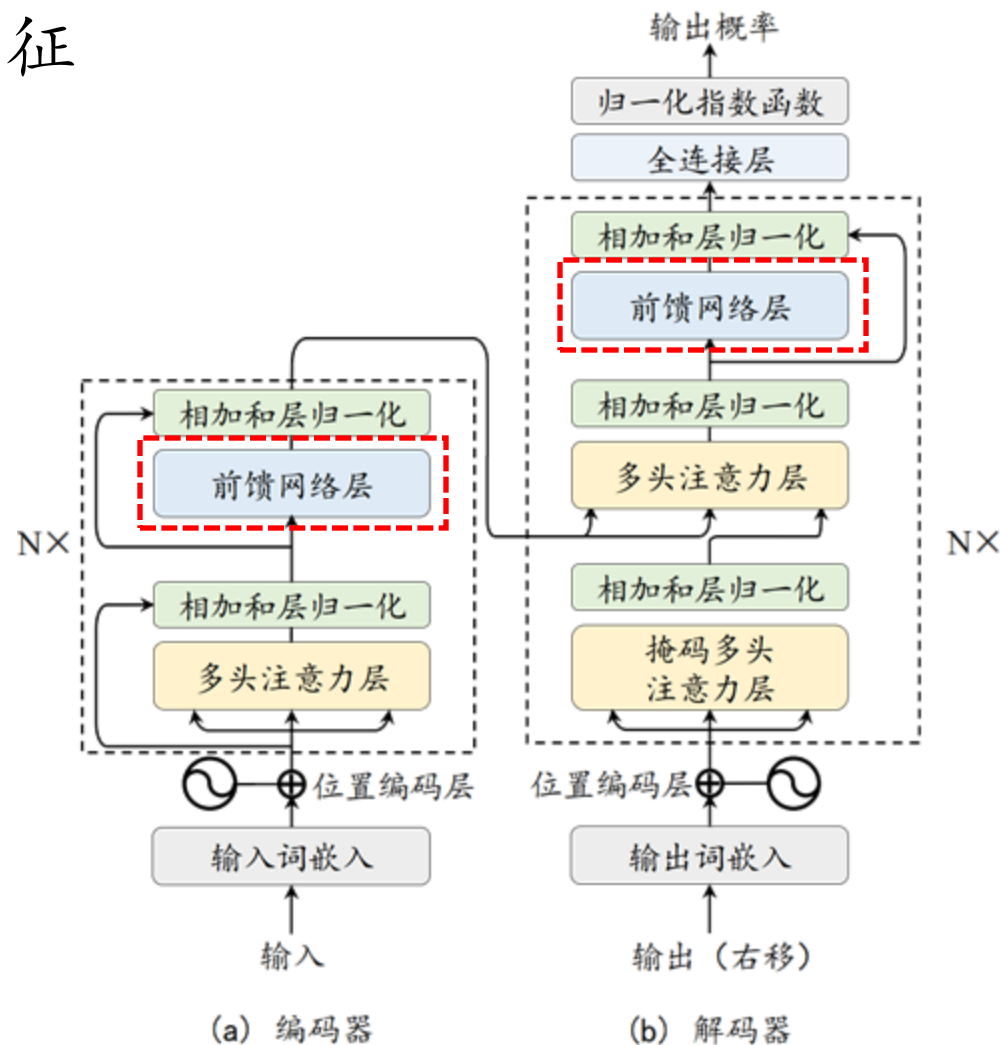
$$\text{FFN}(X) = \sigma(XW^U + b_1)W^D + b_2$$

➤ 线性变换

➤ 先升维、后降维

➤ 非线性激活函数 σ

➤ ReLU 或 GELU 等

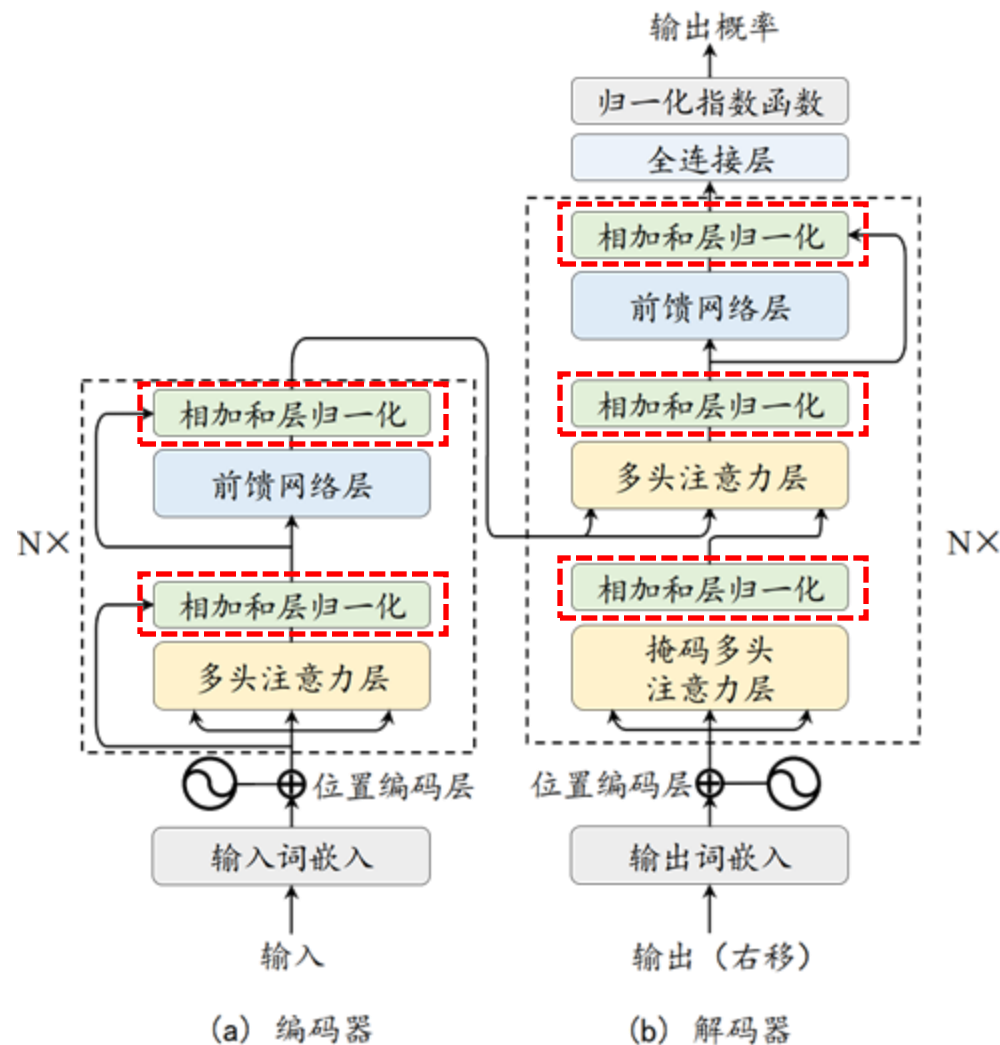


➤ 残差连接

- 将输入与输出相加
- 缓解梯度爆炸和消失

➤ 层归一化

- 对数据进行重新放缩
- 提升训练稳定性



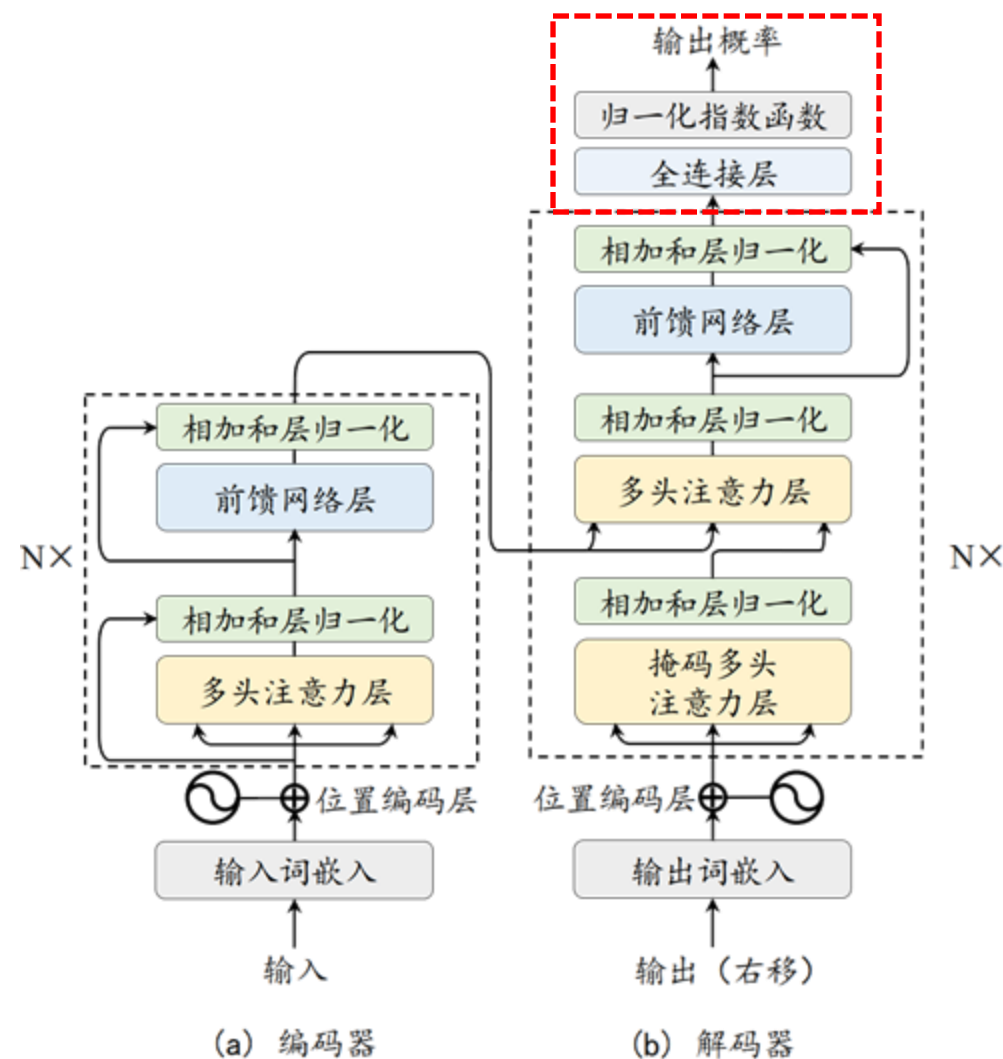
Transformer模型

➤ 输出层：生成词元概率分布

➤ 全连接层

➤ 归一化指数函数 Softmax

$$\mathbf{O} = \text{softmax}(\mathbf{W}^L \mathbf{Y}_L)$$





谢谢