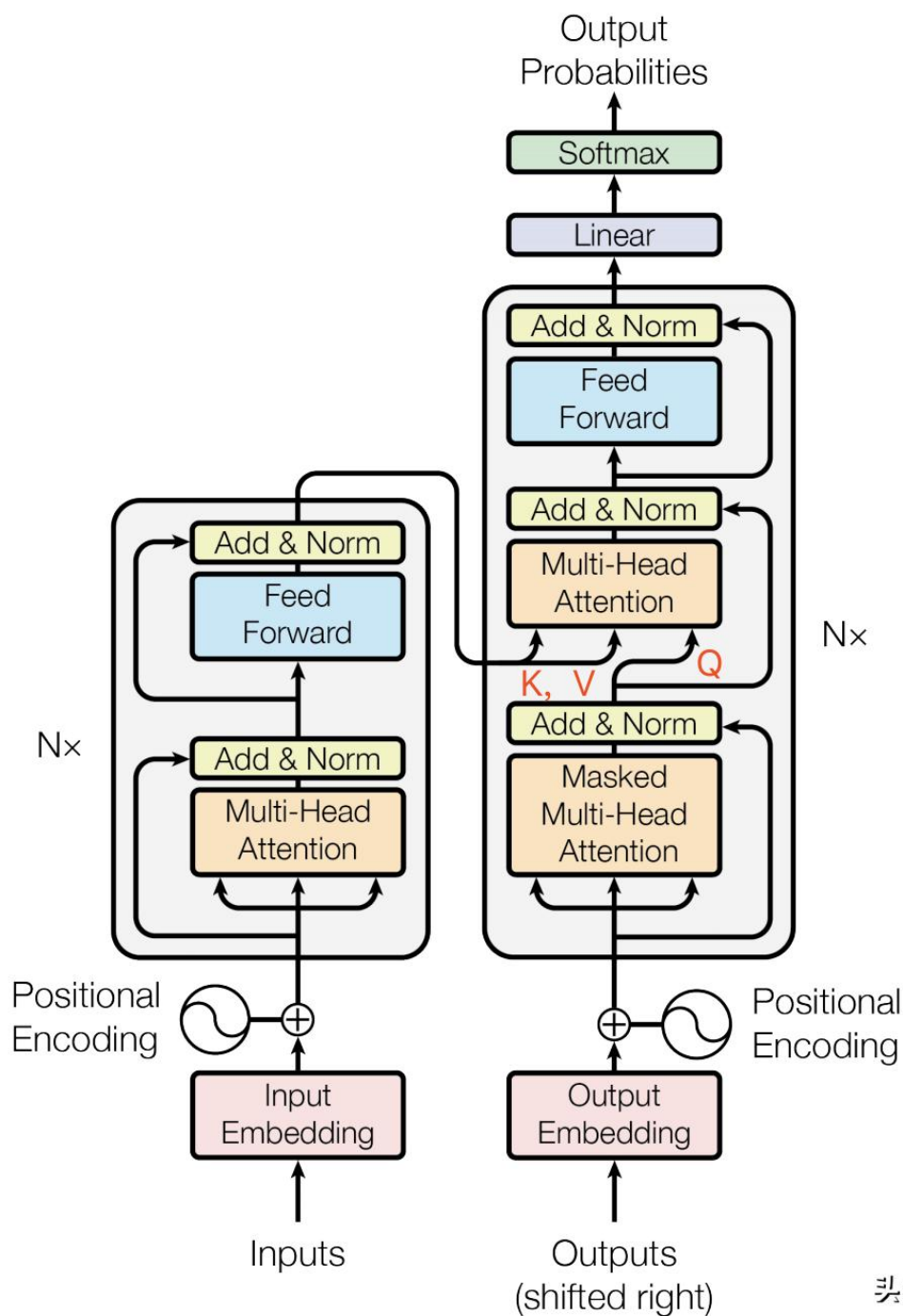


一.Self-Attention自注意力机制



头条 @小狸

1.自注意力机制的基础概念

自注意力机制是一种能够让模型在处理序列数据时，自动关注序列中不同位置之间的依赖关系的机制。它允许序列中的每个元素都与其他元素进行交互，计算出每个元素对其他元素的关注度，从而更好地捕捉序列内部的上下文信息。

例如，在处理句子“他喜欢吃苹果，她喜欢吃香蕉”时，自注意力机制能够让“他”关注到“喜欢吃苹果”，让“她”关注到“喜欢吃香蕉”，同时也能让“苹果”和“香蕉”之间产生一定的关联。

自注意力机制适用于下面第一种N to N

- Each vector has a label.

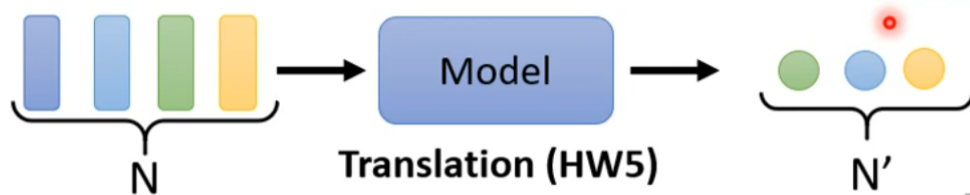


- The whole sequence has a label.



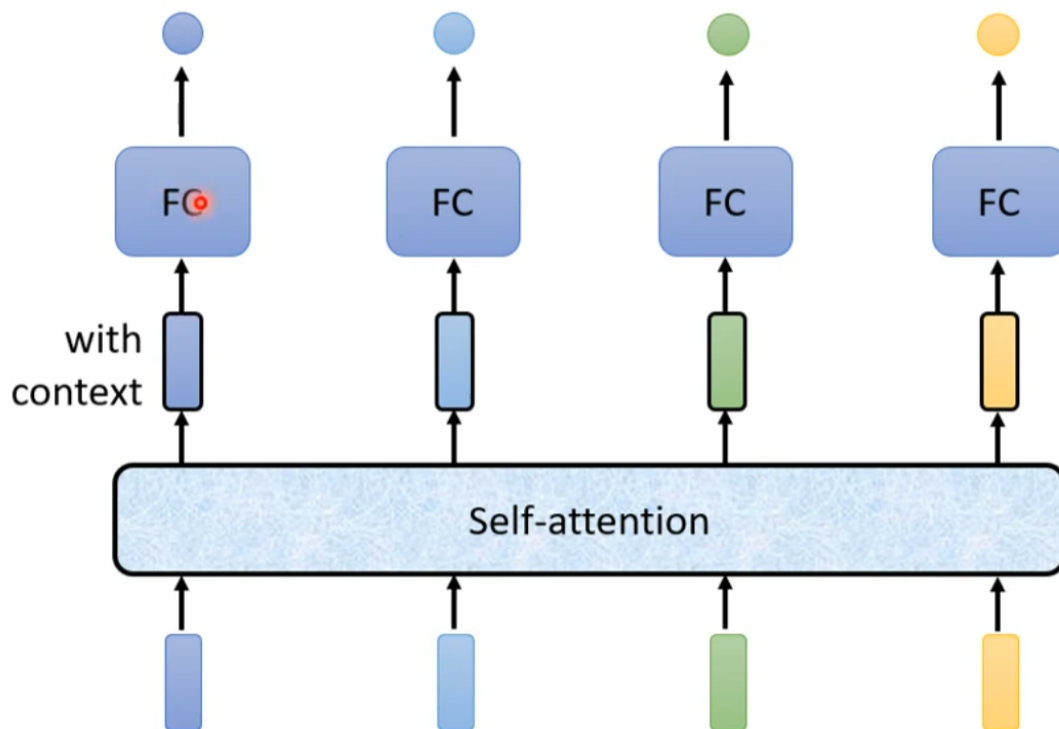
- Model decides the number of labels itself.

seq2seq



9

Self-attention



2. 自注意力机制的核心思想

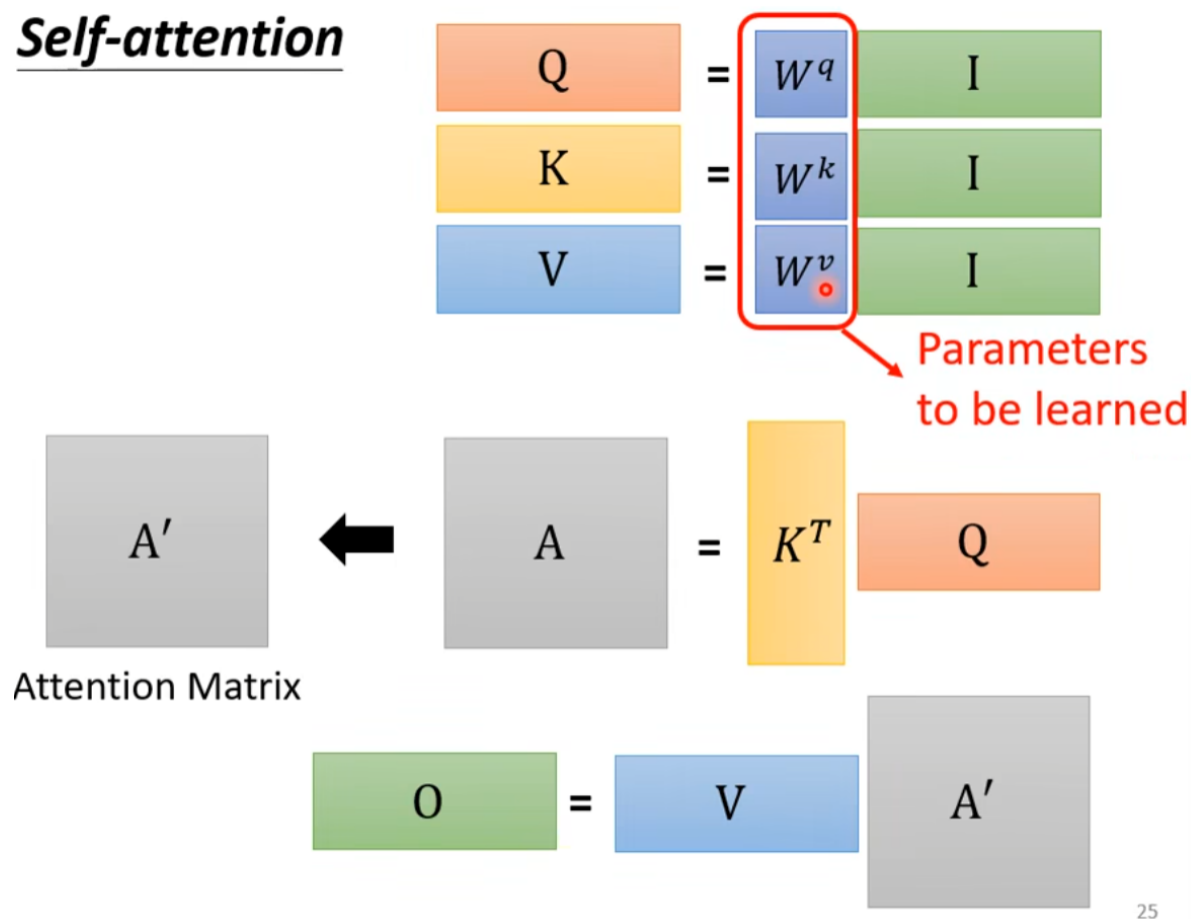
自注意力机制的核心思想是通过计算序列中各个元素之间的相似度（即注意力权重），然后根据这些权重对序列元素进行加权求和，得到每个元素的新表示。这种新表示融合了序列中与该元素相关的其他元素的信息，从而使模型能够更好地理解序列的上下文。

具体来说，对于一个输入序列，自注意力机制会为每个元素生成三个向量：查询向量（Query, Q）、键向量（Key, K）和值向量（Value, V）。然后，通过计算查询向量与键向量的相似度得到注意力权重，再用这些权重对值向量进行加权求和，得到每个元素的输出向量。

3. 自注意力机制的工作原理

总过程用矩阵表示：

Self-attention



25

(一) Q、K、V 的计算

对于输入序列中的每个元素 x_i ，通过三个不同的权重矩阵 W_Q 、 W_K 、 W_V 进行线性变换，分别得到对应的查询向量 q_i 、键向量 k_i 和值向量 v_i ，计算公式如下：
 $q_i = x_i W_Q$ $k_i = x_i W_K$
 $v_i = x_i W_V$

其中， W_Q 、 W_K 、 W_V 是需要通过训练学习得到的参数矩阵。

$$\begin{aligned}
 q^i &= W^q a^i & \begin{matrix} q^1 & q^2 & q^3 & q^4 \\ Q \end{matrix} &= W^q \begin{matrix} a^1 & a^2 & a^3 & a^4 \\ I \end{matrix} \\
 k^i &= W^k a^i & \begin{matrix} k^1 & k^2 & k^3 & k^4 \\ K \end{matrix} &= W^k \begin{matrix} a^1 & a^2 & a^3 & a^4 \\ I \end{matrix} \\
 v^i &= W^v a^i & \begin{matrix} v^1 & v^2 & v^3 & v^4 \\ V \end{matrix} &= W^v \begin{matrix} a^1 & a^2 & a^3 & a^4 \\ I \end{matrix}
 \end{aligned}$$

(二) 注意力权重的计算

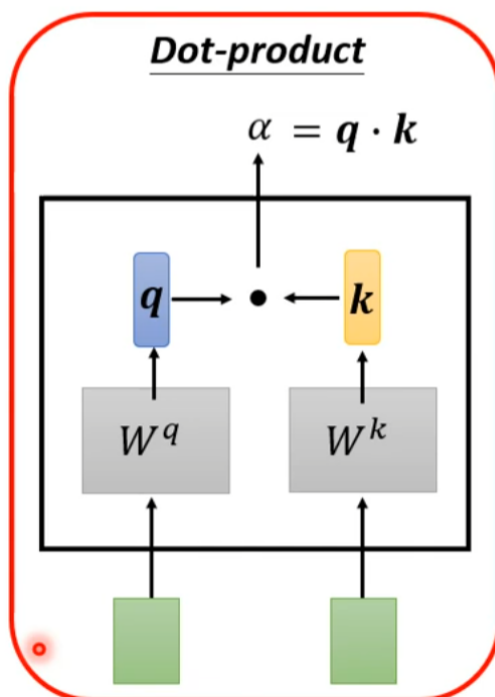
注意力权重用于衡量序列中不同元素之间的关联程度。对于查询向量 q_i ，它与所有键向量 k_j ($j = 1, 2, \dots, n$, n 为序列长度) 的相似度可以通过点积来计算，得到原始的注意力分数：

$$score_{i,j} = q_i \cdot k_j$$

为了使注意力权重在合理的范围内，通常会对原始分数进行缩放，除以键向量维度的平方根 $\sqrt{d_k}$ (d_k 为键向量的维度)，即： $score_{i,j} = \frac{score_{i,j}}{\sqrt{d_k}}$

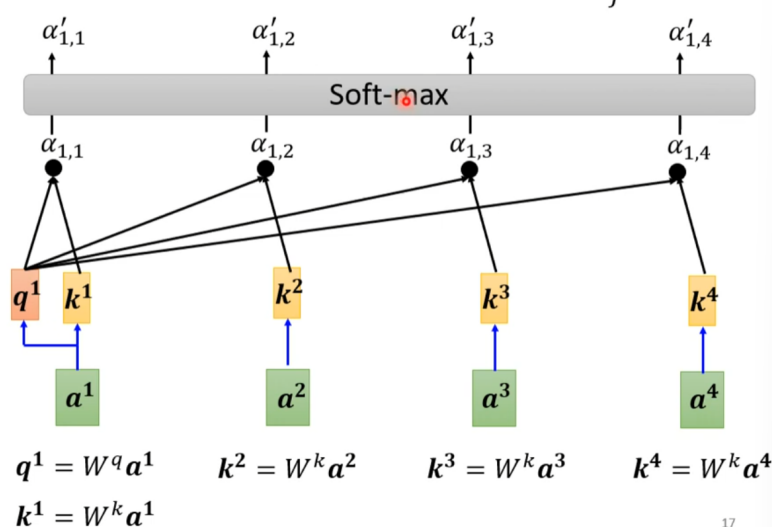
然后，使用 softmax 函数对缩放后的分数进行归一化处理，得到注意力权重 $\alpha_{i,j}$ ，使得权重之和为

$$1: \alpha_{i,j} = \frac{\exp(score_{i,j})}{\sum_{k=1}^n \exp(score_{i,k})}$$



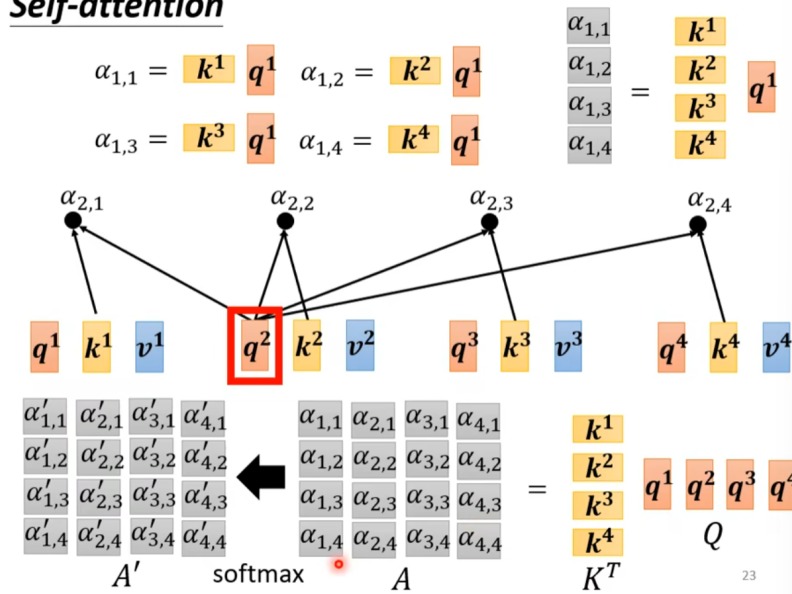
Self-attention

$$\alpha'_{1,i} = \exp(\alpha_{1,i}) / \sum_j \exp(\alpha_{1,j})$$



矩阵乘法表示：

Self-attention

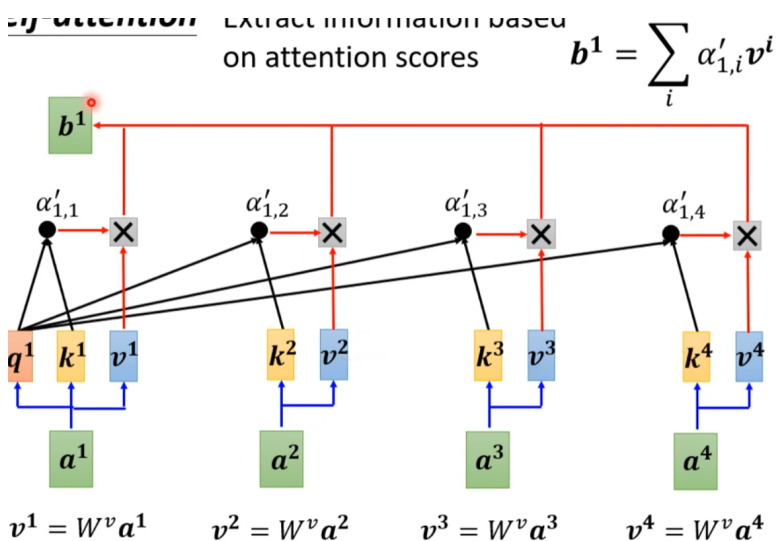


(三) 输出向量的计算

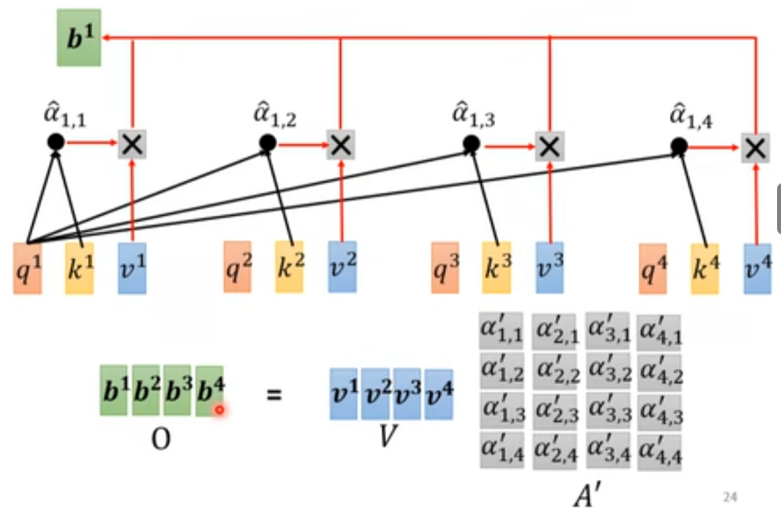
得到注意力权重后，将每个值向量 v_j 按照对应的权重 $\alpha_{i,j}$ 进行加权求和，就可以得到输入元素 x_i

经过自注意力机制处理后的输出向量 o_i ：
$$o_i = \sum_{j=1}^n \alpha_{i,j} v_j$$

对于整个输入序列，将所有输出向量组合起来，就得到了自注意力机制的最终输出矩阵 $O = [o_1, o_2, \dots, o_n]$ 。



矩阵乘法表示：



24

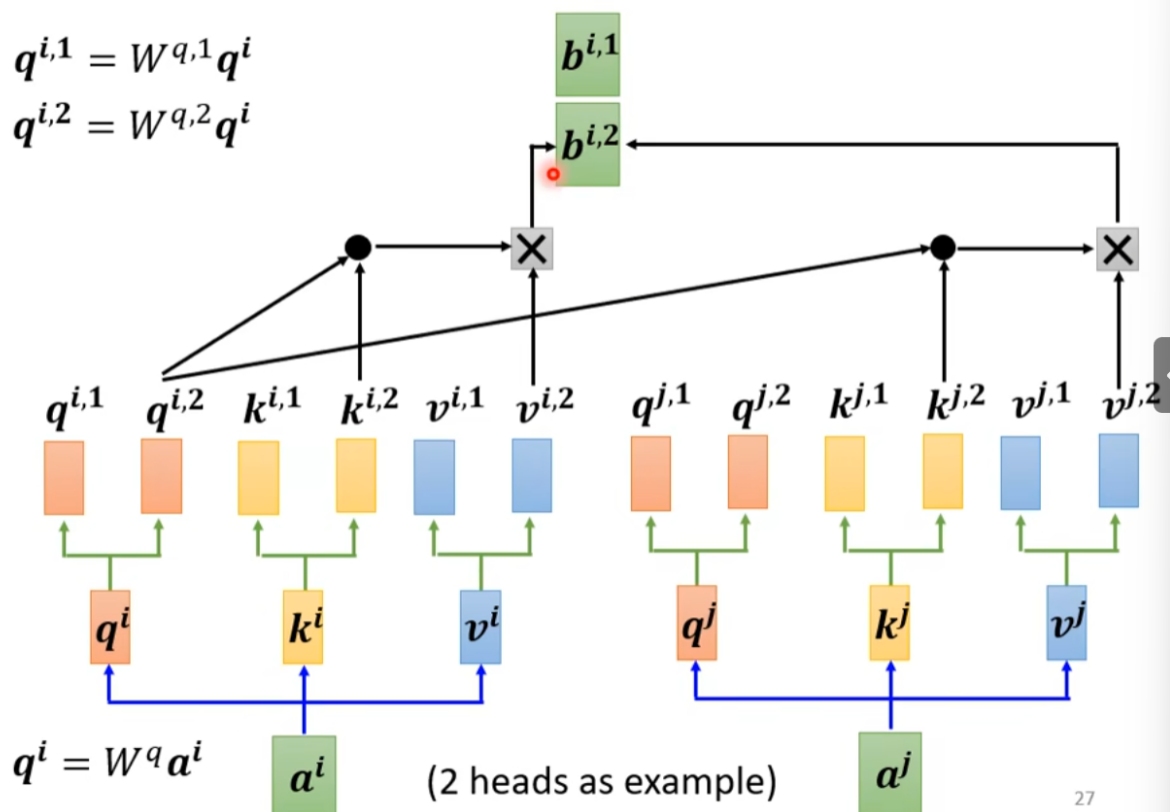
4. 扩充

(一) Multi-head Self-attention

Multi-head Self-attention Different types of relevance

$$q^{i,1} = W^{q,1} q^i$$

$$q^{i,2} = W^{q,2} q^i$$



27

$$b^i = W^O \begin{bmatrix} b^{i,1} \\ b^{i,2} \end{bmatrix}$$

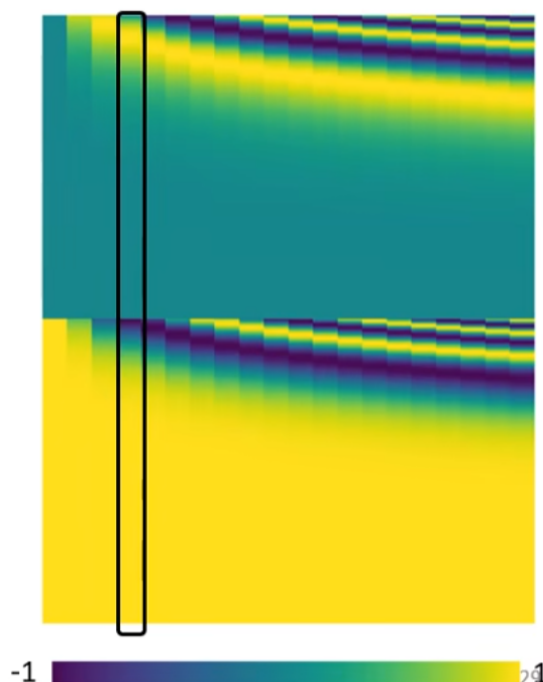
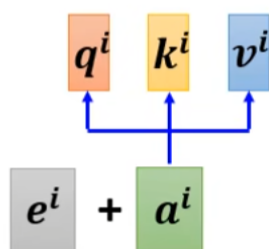
(二) 位置编码Positional Encoding

在自注意力机制中，模型本身并不能直接捕捉序列元素的位置信息，因为自注意力的计算是对序列中所有元素同时进行的，不依赖于元素的顺序。而对于很多序列数据（如自然语言句子），元素的位置信息至关重要，不同的位置会导致序列的含义发生变化。例如，“我喜欢你”和“你喜欢我”，词语相同但位置不同，含义完全相反。因此，需要通过位置编码为序列中的每个元素添加位置信息，让模型能够区分不同位置的元素。

Positional Encoding

Each column represents a positional vector e^i

- No position information in self-attention.
- Each position has a unique positional vector e^i
- **hand-crafted**



5.自注意力机制与传统注意力机制的对比

(一) 传统注意力机制

传统的注意力机制通常是在两个不同的序列之间计算注意力，例如在机器翻译中，解码器在生成每个词时，会关注编码器输出的源语言序列中的相关词。其注意力计算是跨序列的，需要一个查询序列和一个键值序列。

(二) 自注意力机制

自注意力机制则是在同一个序列内部计算注意力，查询序列、键序列和值序列都来自同一个输入序列。这种机制能够更好地捕捉序列内部的长距离依赖关系，并且可以实现并行计算，大大提高了模型的训练和推理效率。

相比之下，自注意力机制在处理长序列时表现更优，因为它不需要像 RNN 那样按顺序处理序列，可以同时处理序列中的所有元素。同时，它也避免了卷积神经网络（CNN）中感受野有限的问题，能够捕捉到序列中任意位置之间的依赖关系。

6.自注意力机制的应用

自注意力机制在自然语言处理领域得到了广泛的应用，除了 Transformer 模型外，在文本分类、命名实体识别、情感分析等任务中都有出色的表现。

在计算机视觉领域，自注意力机制也被用于图像分类、目标检测等任务，通过捕捉图像中不同区域之间的依赖关系，提高模型对图像的理解能力。

二.类神经网络训练不起来怎么办

1.先搞懂最基础的：Encoder-Decoder 框架

其实啊，Encoder-Decoder（编码器 - 解码器）不是一个具体模型，而是一种**通用的“设计思路”**，就像盖房子的“框架”—— 你可以用砖（RNN）、用混凝土（Transformer）来填，但整体结构不变。

核心思想：“理解”再“生成”

想象一个场景：把中文“我爱学习”翻译成英文“I love studying”。这个过程可以拆成两步：

- 第一步：先“读懂”中文这句话的意思（比如理解“我”是主语，“爱”是谓语，“学习”是宾语）；
- 第二步：根据这个“意思”，生成对应的英文句子。

Encoder-Decoder 框架就是干这个的：

- **Encoder（编码器）**：负责“读懂”输入（比如中文句子），把输入转换成一个“有意义的向量”（叫“上下文向量”，context vector）；
- **Decoder（解码器）**：负责根据这个“上下文向量”，生成输出（比如英文句子）。

直观图示（简化版）

输入序列（如中文） → Encoder → 上下文向量（context） → Decoder → 输出序列（如英文）

适用场景

只要是“输入和输出都是序列”的任务，几乎都能用这个框架：

- 机器翻译（输入一种语言，输出另一种语言）；
- 文本摘要（输入长文本，输出短摘要）；
- 问答系统（输入问题，输出答案）；
- 语音识别（输入语音波形，输出文字）。

简单说：只要需要“先理解输入，再生成输出”，Encoder-Decoder 就能上。

2.Seq2Seq 模型：Encoder-Decoder 的“初代网红”实现

有了 Encoder-Decoder 框架，总得有具体的模型来实现它吧？Seq2Seq（Sequence to Sequence）就是最经典的一个，2014 年由谷歌团队提出，当年在机器翻译上效果惊艳。

Seq2Seq 的“内核”：用 RNN 家族填充框架

Seq2Seq 的 Encoder 和 Decoder，用的是**RNN（循环神经网络）** 或者它的升级版 LSTM、GRU（因为 RNN 容易“忘事”，LSTM/GRU 能记住长一点的信息）。

具体流程：

- Encoder（用 LSTM 举例）：
 1. 输入序列（比如“我 / 爱 / 学习”）每个词依次进 LSTM；
 2. 每一步 LSTM 会输出一个“隐藏状态”，但 Encoder 只保留**最后一个隐藏状态**作为“上下文向量”（代表对整个输入的理解）。
- Decoder（也用 LSTM）：
 1. 以 Encoder 输出的“上下文向量”作为初始隐藏状态；
 2. 一步步生成输出序列：先输出第一个词（比如“l”），再把“l”作为输入，生成下一个词（“love”），直到生成结束符（比如“”）。这个过程叫“自回归生成”。

Seq2Seq 的“痛点”

虽然当年很牛，但用久了发现问题：

- **“信息瓶颈”**：Encoder 把所有输入信息压缩成一个固定长度的向量（比如 100 维），如果输入很长（比如一篇文章），这个向量根本装不下所有信息，后面的词会被“忘光”；
- **“并行性差”**：RNN 是“顺序计算”的（必须等前一个词处理完才能处理下一个），没法像 CNN 那样并行计算，训练慢得很。

举个例子：如果输入是“从前有座山，山里有座庙，庙里有个老和尚...（1000 字故事）”，Seq2Seq 的 Encoder 最后输出的向量，可能只记得“老和尚”，前面的“山”和“庙”早忘了。

3.Transformer：Encoder-Decoder 的“终极进化”

2017 年，谷歌又放了个大招——Transformer 模型（论文《Attention Is All You Need》），直接抛弃了 RNN，改用**自注意力机制（Self-Attention）**，解决了 Seq2Seq 的痛点，现在大模型（比如 GPT、BERT）几乎都基于它。

为什么 Transformer 能“封神”？

核心原因：用**“注意力”**代替 RNN，解决了并行和长序列问题。

想象一下：读句子时，你不会平均分配注意力，比如“猫追狗”，你会重点关注“猫”和“狗”的关系；Transformer 的“注意力”就是干这个的——让模型自动学会“该关注输入中的哪些词”。

Transformer 的整体结构（也是 Encoder-Decoder 框架！）

Transformer 依然遵循 Encoder-Decoder 框架，但内部细节大改：

输入序列 → Encoder 部分（6个Encoder层堆叠） → 上下文向量 → Decoder 部分（6个Decoder层堆叠） → 输出序列

(1) 先看输入：词嵌入 + 位置编码

不管 Encoder 还是 Decoder，输入的词都要先处理成向量：

- **词嵌入 (Word Embedding)**：把每个词转换成固定长度的向量（比如“猫”→ [0.2, 0.5, ..., 0.1]），这一步和 Seq2Seq 一样；
- **位置编码 (Positional Encoding)**：RNN 是按顺序处理的，天然知道词的位置（第一个词、第二个词），但 Transformer 并行计算，不知道顺序！所以必须手动加“位置信息”，比如用正弦余弦函数生成位置向量，和词嵌入加起来。

(2) Encoder 层：“吃透”输入序列

每个 Encoder 层有两个核心模块（加了残差连接和层归一化）：

- 多头自注意力 (Multi-Head Self-Attention)：
 - 作用：让每个词“关注输入序列中其他相关的词”。比如“他喜欢篮球，每天都打它”，“它”要关注“篮球”；
 - “多头”：把注意力分成 8 组（论文里是 8 头），每组学不同的“关注角度”（比如一组关注语法，一组关注语义），最后拼接起来，信息更全面；
- **前馈神经网络 (Feed Forward Network)**：对每个词的向量做一次非线性转换（比如先升维再降维），增强模型能力。

小提示：Encoder 的自注意力是“双向的”——每个词能看到所有词（包括前面和后面的），所以能“吃透”整个输入的上下文。

(3) Decoder 层：“精准”生成输出

每个 Decoder 层有三个核心模块：

- 掩蔽多头自注意力 (Masked Multi-Head Self-Attention)：
 - 作用：生成输出时，只能“关注已经生成的词”（比如生成第 3 个词时，只能看第 1、2 个词），不能“偷看”后面的词（否则就作弊了）；
 - “掩蔽”就是用一个矩阵把“未来的词”遮住，让模型看不到；
- 编码器 - 解码器注意力 (Encoder-Decoder Attention)：
 - 作用：让 Decoder 生成每个词时，关注 Encoder 输出的“输入序列中相关的词”。比如翻译时，生成英文“I”要关注中文“我”；
- **前馈神经网络**：和 Encoder 里的一样，做非线性转换。

3. Transformer 的优势 (对比 Seq2Seq)

特点	Seq2Seq (RNN)	Transformer (自注意力)
并行性	差 (必须按顺序算)	好 (所有词可以同时处理)
长序列处理	差 (固定向量存不下信息)	好 (注意力能关注任意位置)
上下文理解	局部 (主要记最近的词)	全局 (能看到所有词的关系)

4.总结：三者的关系

1. **Encoder-Decoder 是“设计框架”**：定义了“先编码输入，再解码输出”的流程；
2. **Seq2Seq 是这个框架的“早期实现”**：用 RNN/LSTM 填充，解决了很多序列转换问题，但有瓶颈；

3. **Transformer 是“更牛的实现”**：用自注意力代替 RNN，解决了 Seq2Seq 的痛点，成为现在大模型的基础。

简单说：Encoder-Decoder 是“骨架”，Seq2Seq 和 Transformer 是“不同的血肉”，后者更健壮