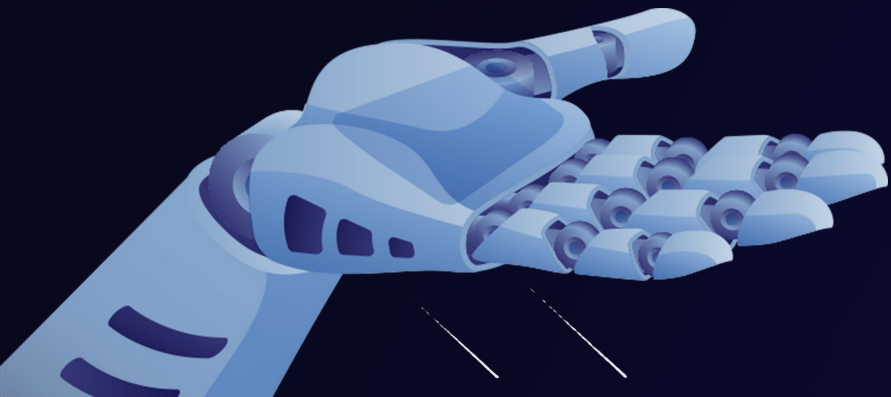


# 第四章：Transformer详解及翻译模型搭建





# CONTENTS



4.1 Transformer简介：基于自注意力机制的Seq2Seq模型

4.2 自注意力机制：充分考虑上下文

4.3 Transformer模型架构：输入、编码器、解码器和输出

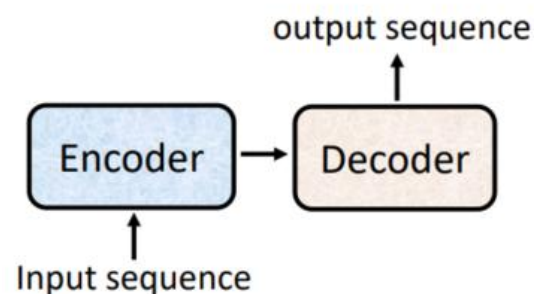
4.4 Transformer训练过程

4.5 Transformer代码实操：了解大模型奠基性算法内核

## 4.1 Transformer简介：基于自注意力机制的Seq2Seq模型

### 编码器-解码器

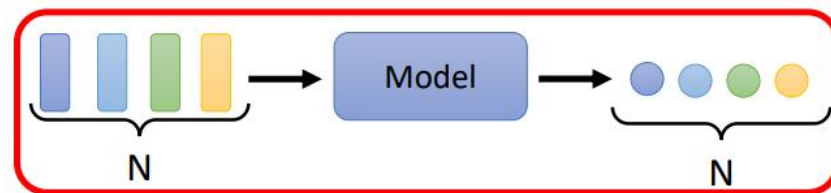
编码器-解码器架构由两部分组成：编码器和解码器。编码器负责将输入序列编码为一个固定长度的向量，而解码器则负责从这个向量中解码出输出序列。这种架构对于处理可变长度的输入和输出序列非常有效。



### 序列传导模型

- 1) 输入N个词，输出N个词。  
例如：词性标注、NER。
- 2) 输入N个词，输出一个词。  
例如：句子分类。
- 3) 输入N个词，输出N'个词。  
例如：翻译、语音识别、语音合成。

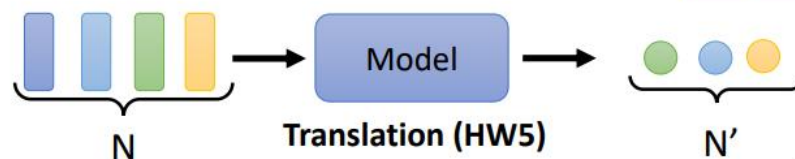
- Each vector has a label. focus of this lecture



- The whole sequence has a label.



- Model decides the number of labels itself. seq2seq



## 4.1 Transformer简介：基于自注意力机制的Seq2Seq模型

### RNN

由于循环神经网络的顺序性质，  
训练不能并行化。



### CNN

可以并行化，但是学习远距离  
位置之间的依赖关系比较困难



LSTM、GRU。。。

ByteNet、ConvS2S。。。

### Transformer

- 1) 训练并行化，且可以学习到远距离依赖。
- 2) 第一个完全依靠**自注意力机制**计算输入和输出表示而不使用序列对齐RNN或卷积的**序列传导模型**。

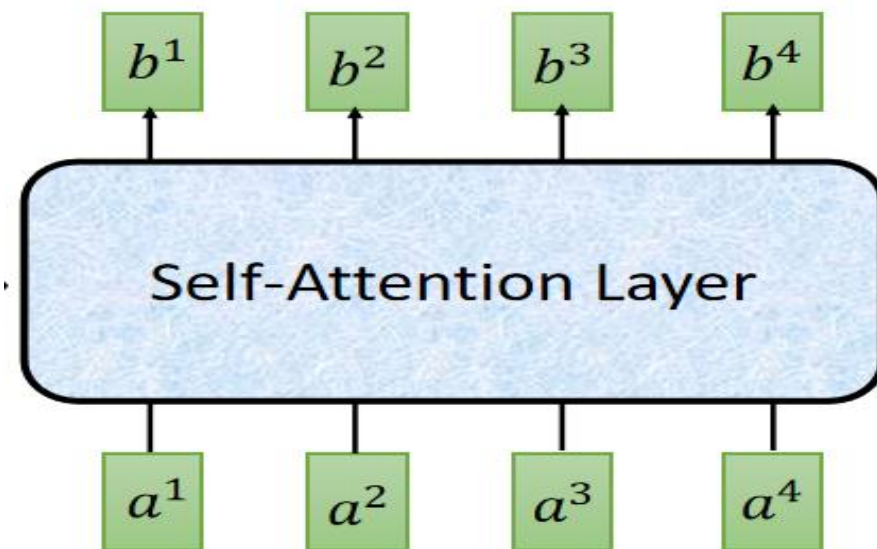
## 4.2 自注意力机制：充分考虑上下文 【简介】

自注意力机制主要用于处理序列数据，如自然语言处理中的文本。

自注意力机制通过计算每个词与其他词之间的相似度来建立它们之间的关系，并根据这些关系来加权地计算每个词的表示。

优势一：考虑上下文。

优势二：计算速度快。



## 4.2 自注意力机制：充分考虑上下文

### 【Q K V计算】

1) 原始输入 $x^i$ 为一个词或token。

2)  $x^i$ 经过转换，变成新的向量 $a^i$ 。

3)  $a^i$ 通过转换，计算得到 $q^i$ 、 $k^i$ 和 $v^i$ 。

$q$ : query (to match others)

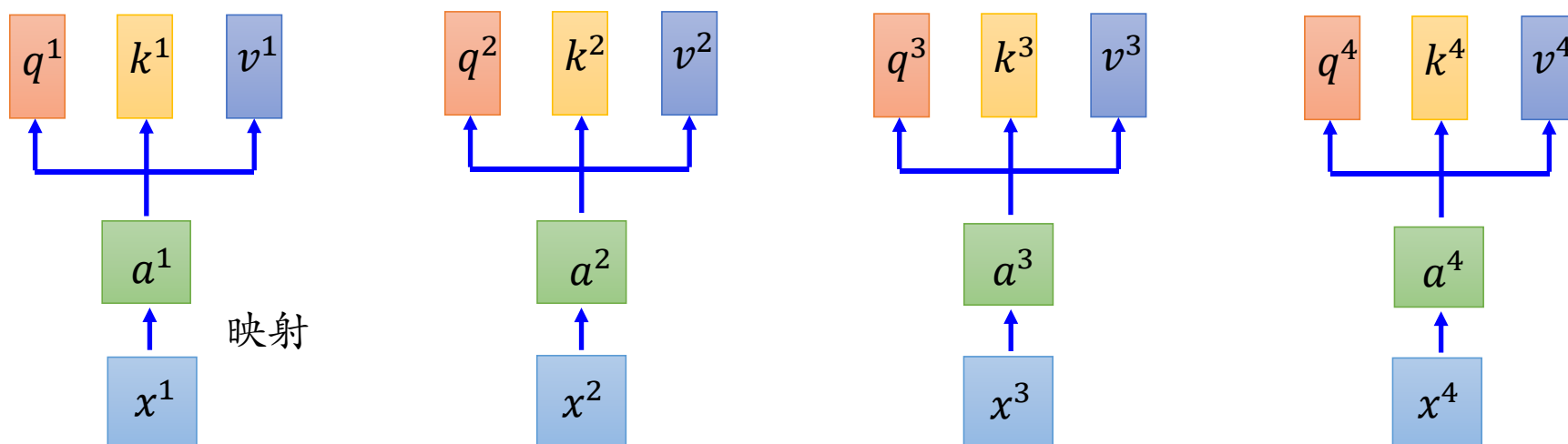
$$q^i = W^q a^i$$

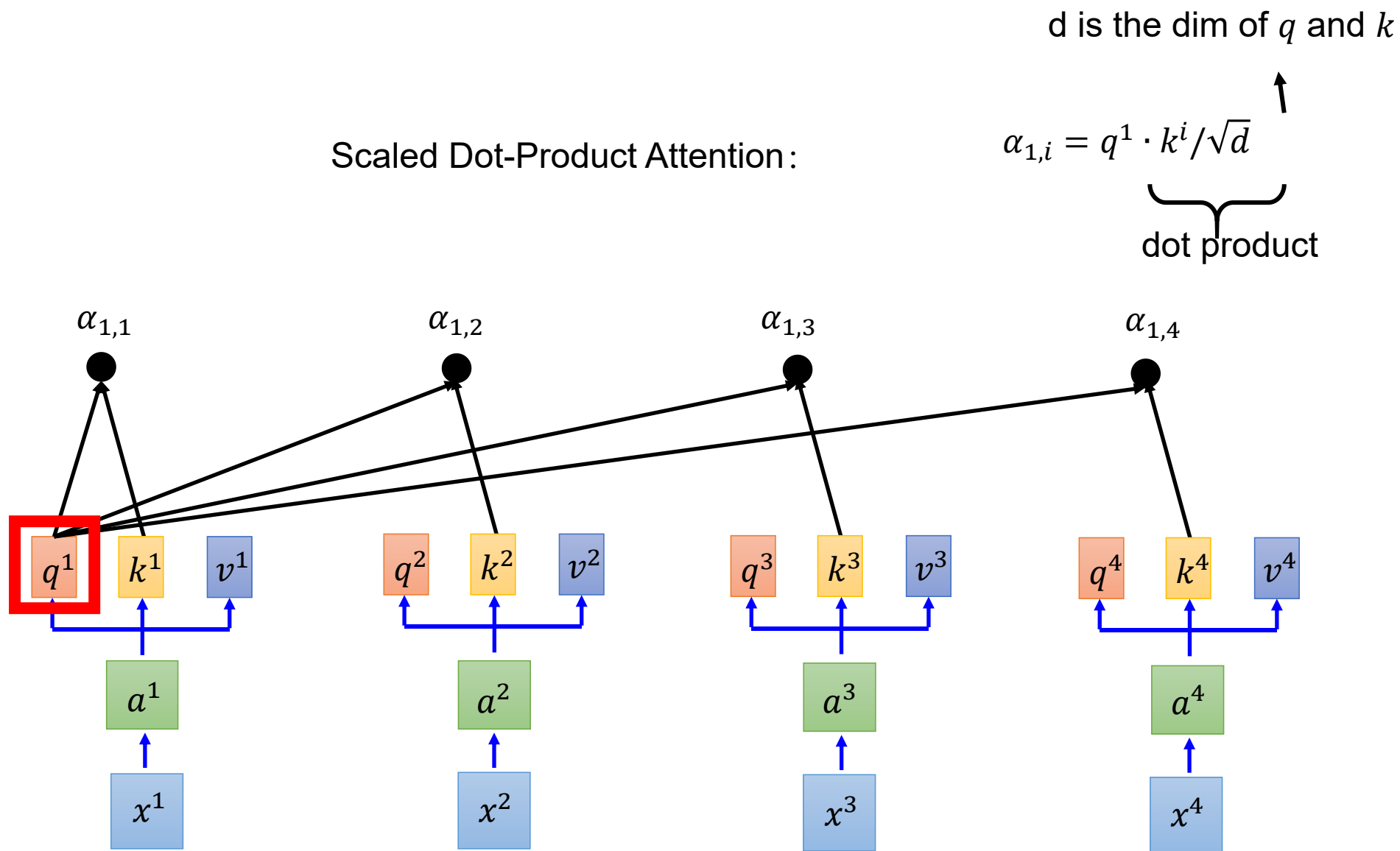
$k$ : key (to be matched)

$$k^i = W^k a^i$$

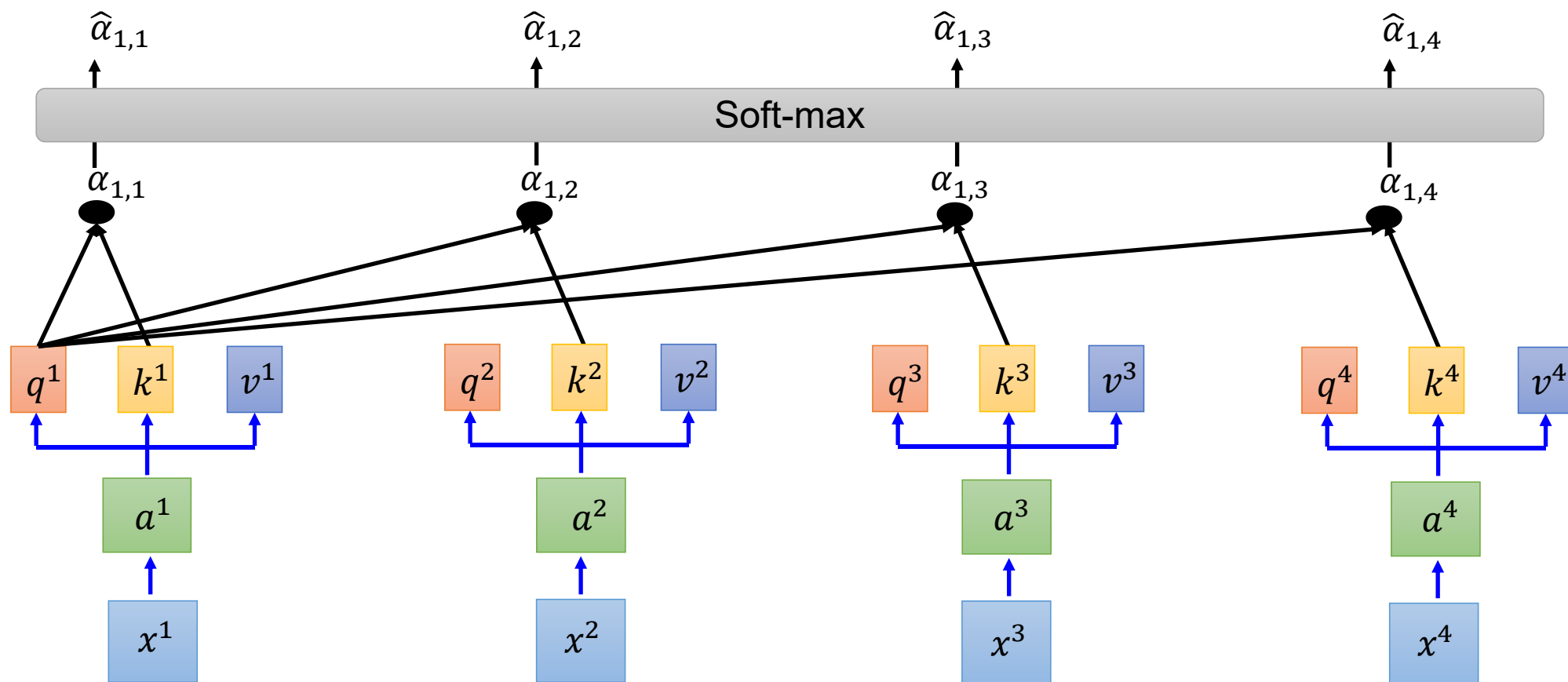
$v$ : information to be extracted

$$v^i = W^v a^i$$

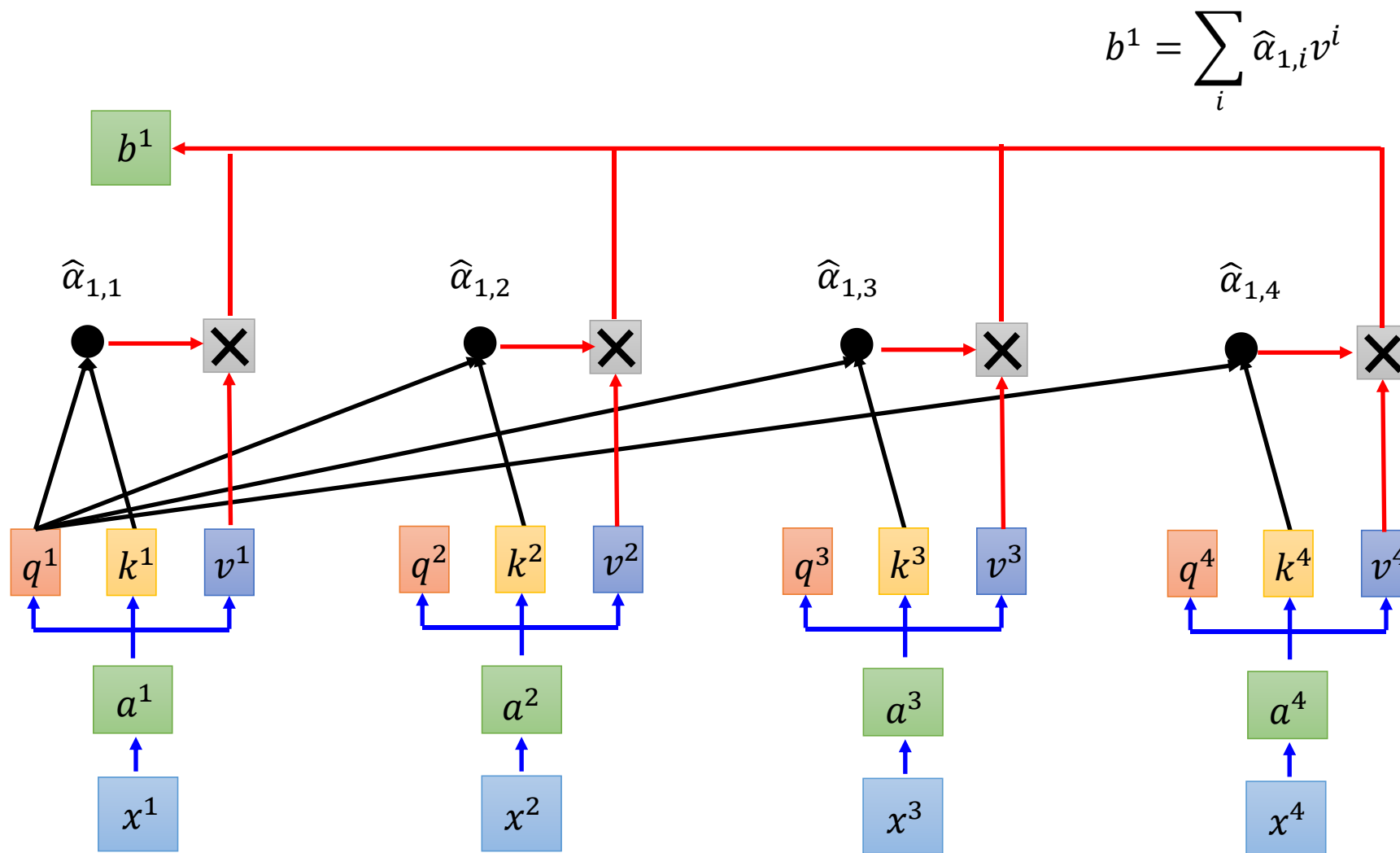


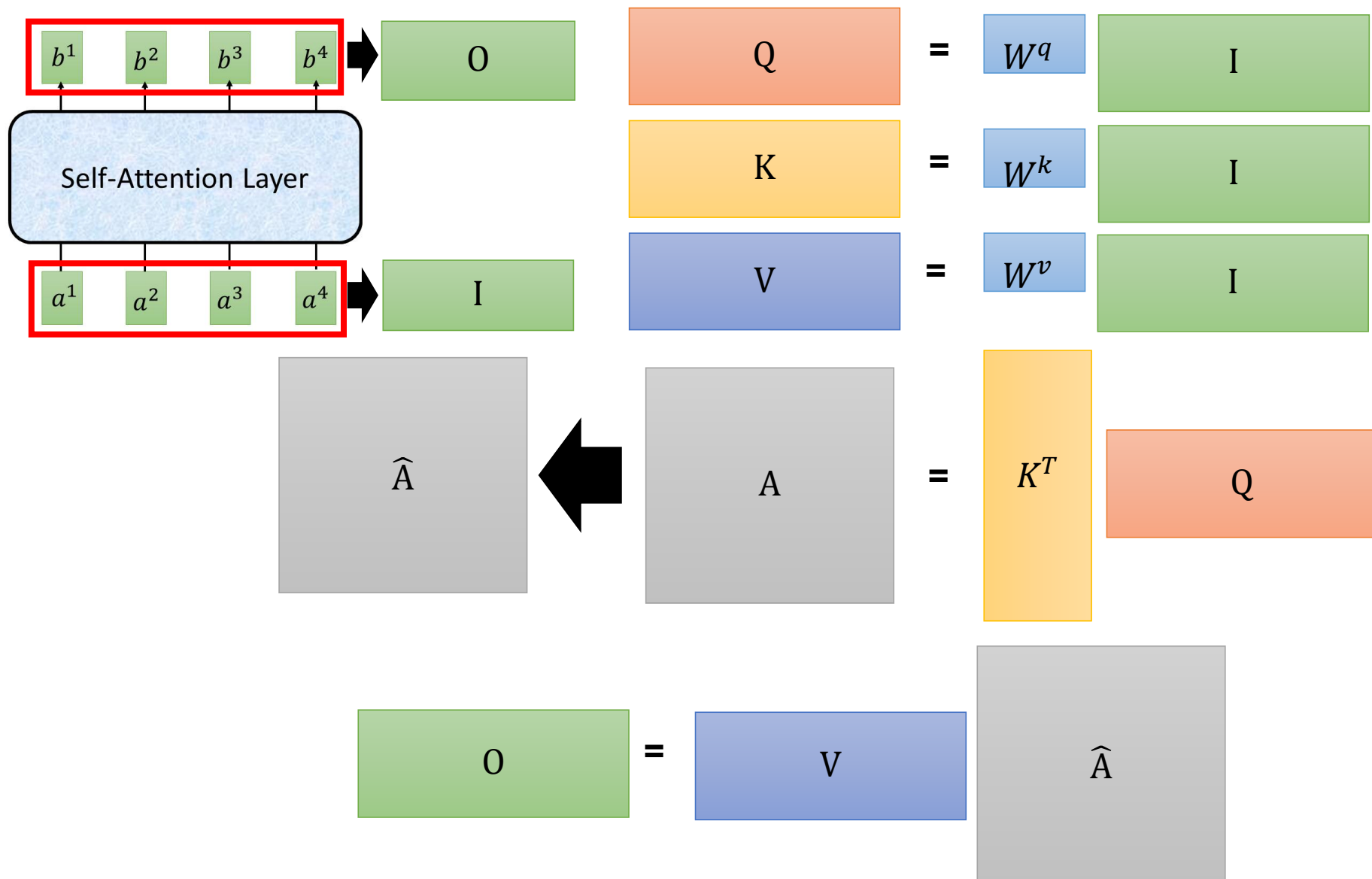


$$\hat{\alpha}_{1,i} = \exp(\alpha_{1,i}) / \sum_j \exp(\alpha_{1,j})$$



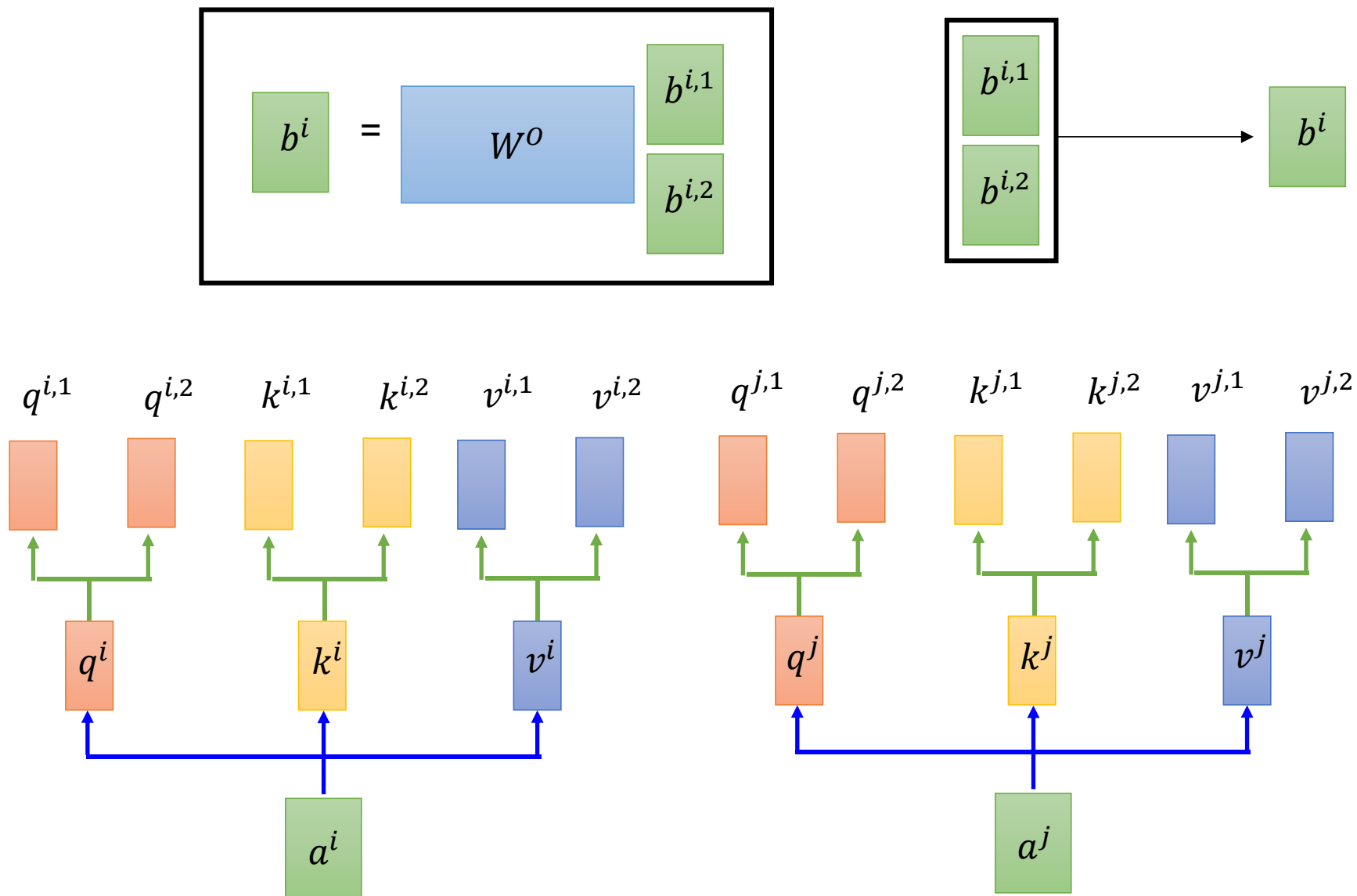






## 4.2 自注意力机制：充分考虑上下文

【多头】



## 4.3 Transformer模型架构 【概况】

左侧：

1) 原始输入部分

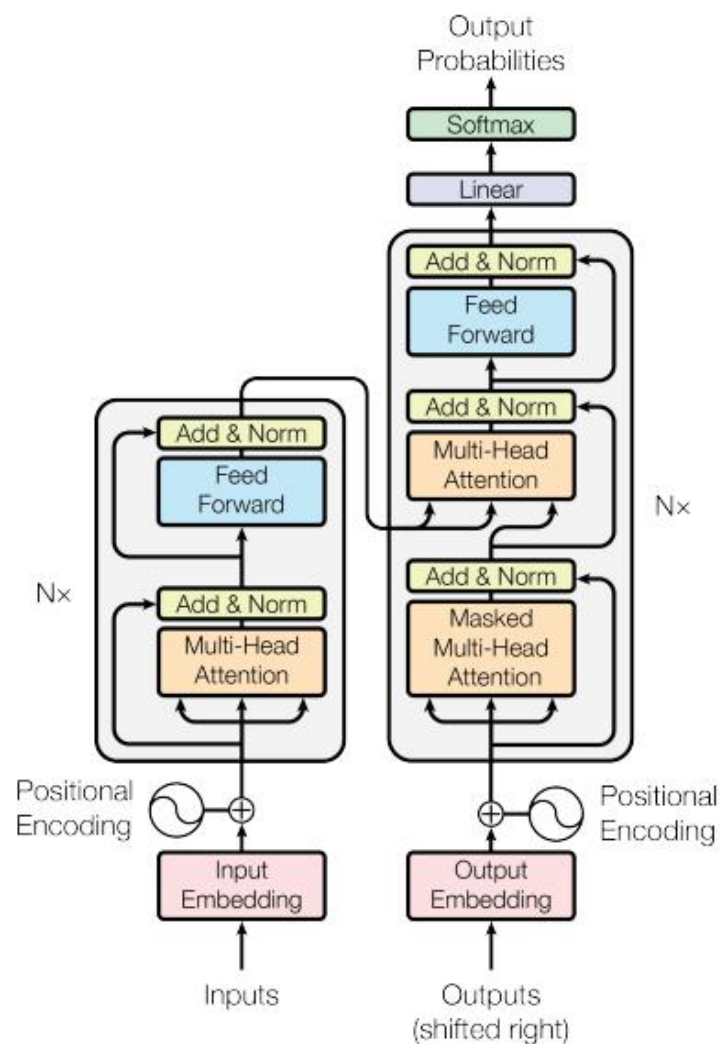
2) 编码器部分

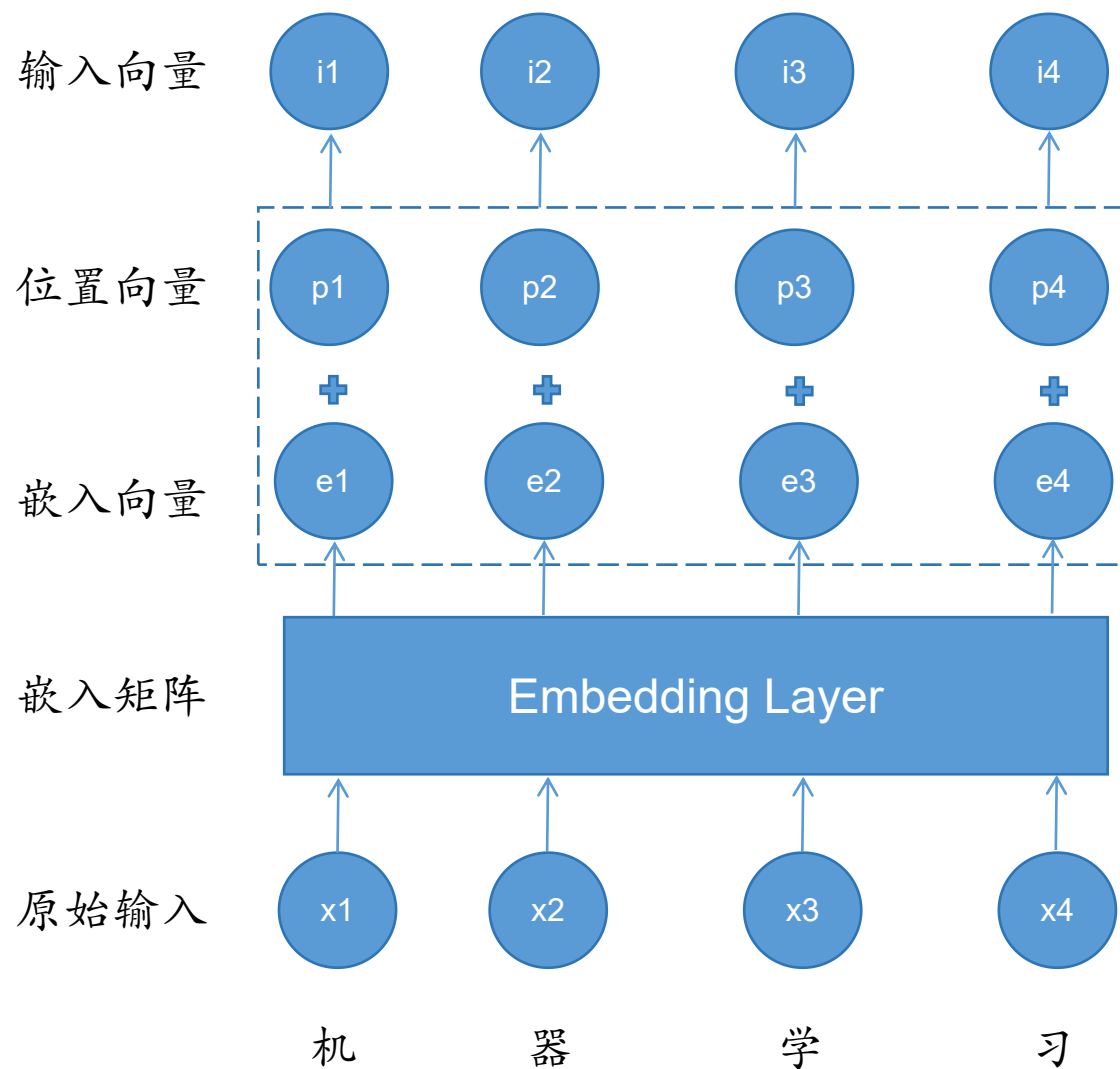
右侧：

3) 目标输入部分

4) 解码器部分

5) 目标输出部分



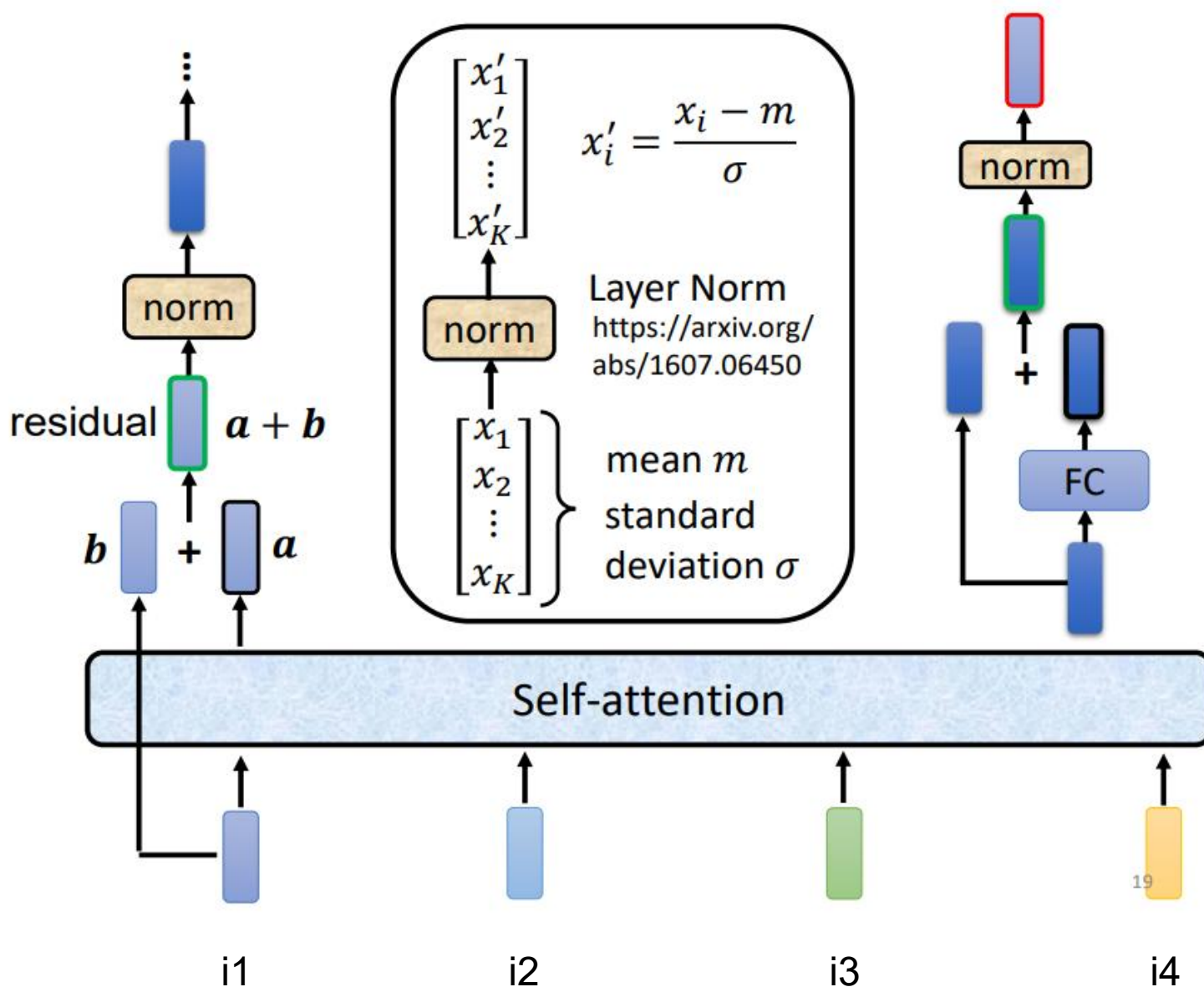


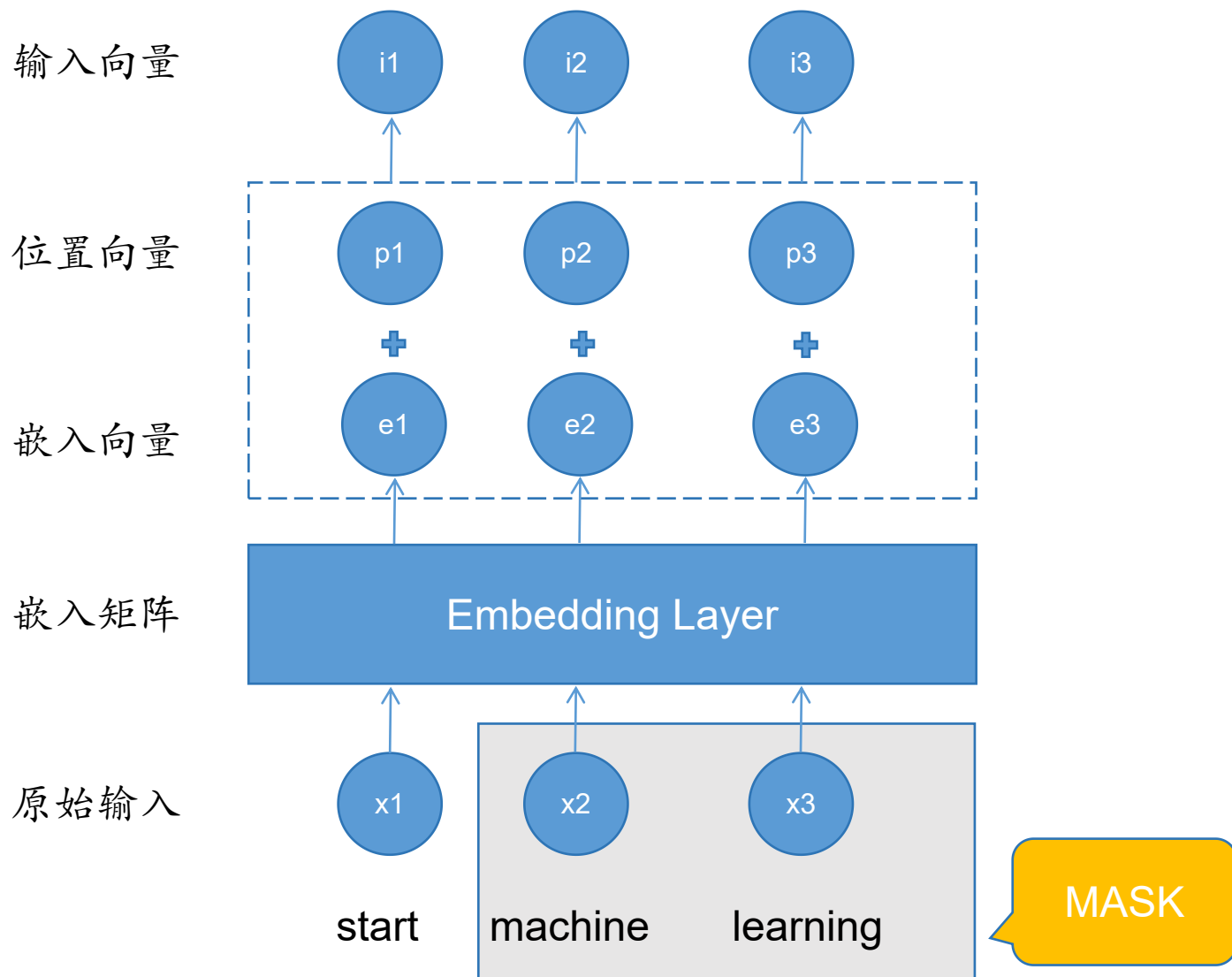
### 位置向量

- 1) 位置编码具有与嵌入相同的 $d_{\text{model}}$ 维度，以便两者可以相加。
- 2) Transformer使用不同波长的正弦和余弦函数来作为位置编码。
- 3) 公式如下， $\text{pos}$ 是位置， $i$ 是维度。也就是说，位置编码的每个维度对应一个正弦或余弦曲线。波长形成了一个几何级数，从 $2\pi$ 到 $10000\pi$ 。
- 4) 容易学习到通过相对位置进行注意力计算，因为对于任何固定的偏移量 $k$ ， $\text{PE}(\text{pos}+k)$ 可以被表示为 $\text{PE}(\text{pos})$ 的线性函数。
- 5) 允许模型外推到比训练中遇到的序列长度更长的序列长度。

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

- 1) 输入向量首先进行自注意力计算。
- 2) 自注意力输出的向量与原始向量进行残差计算。
- 3) 残差计算输出的向量进行层标准化。
- 4) 层标准化输出的向量进行全连接计算。
- 5) 全连接计算输出的向量进行残差计算。
- 6) 第5) 步的残差计算后再进行层标准化计算。





与原始输入部分基本相同，区别在于需要增加自回归当前词信息和掩码信息。

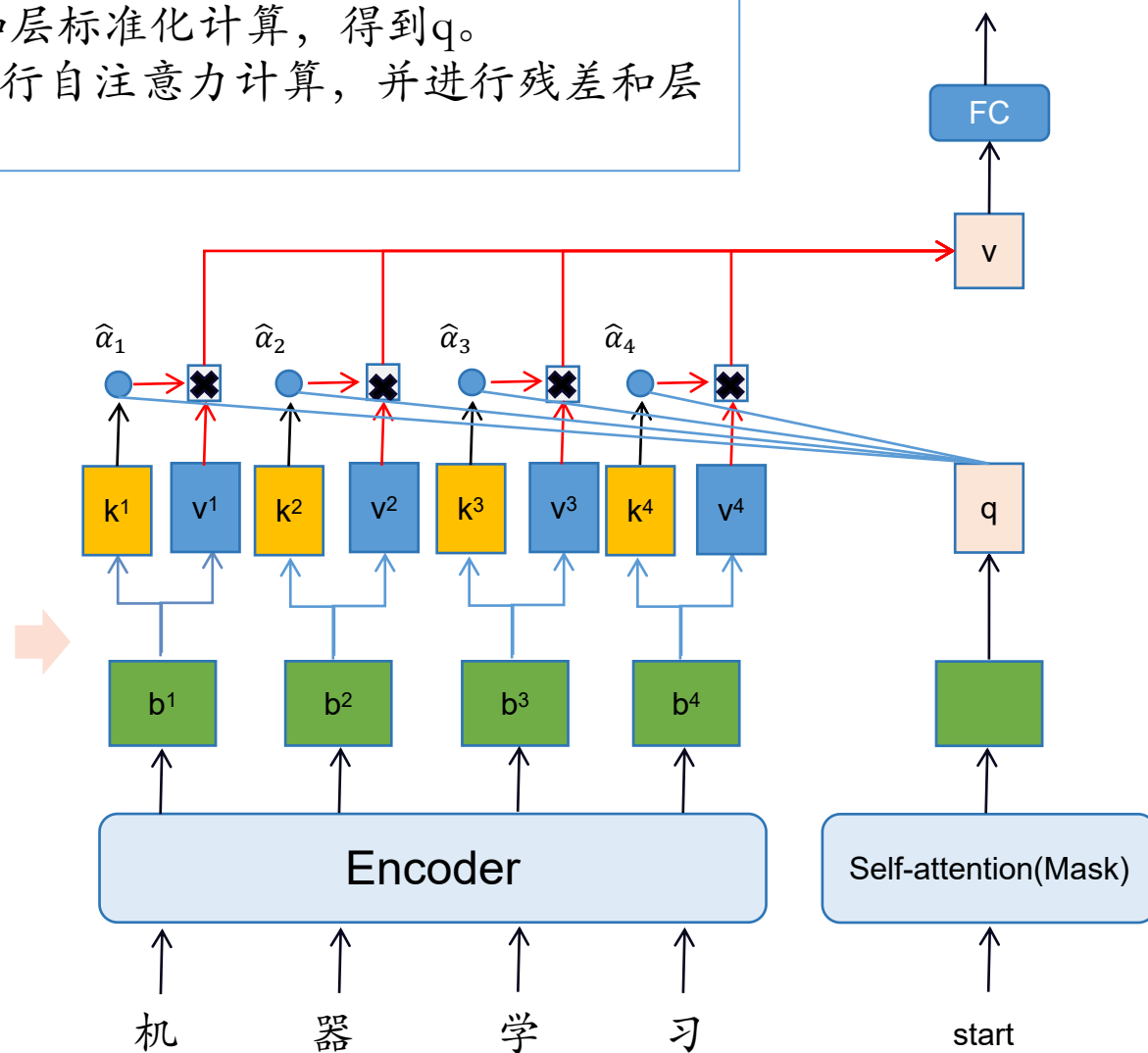
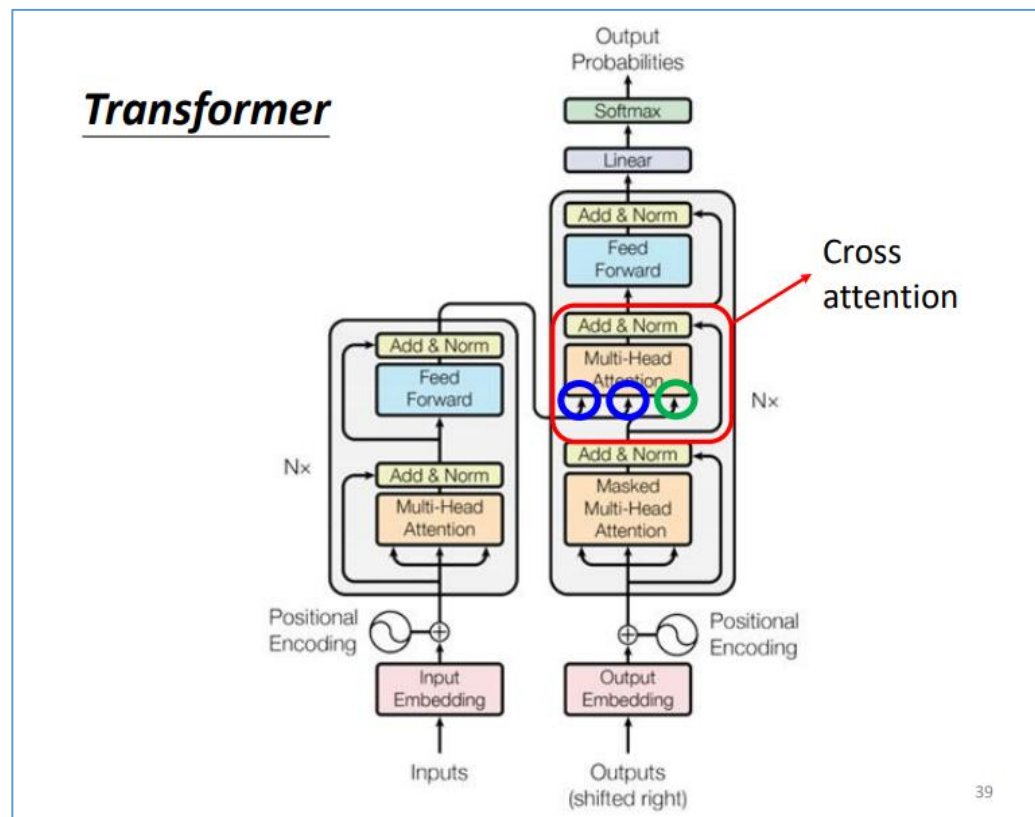
自回归 (Autoregressive) 是指解码器在生成每个输出时，依赖于之前已经生成的输出。这种方法常用于自然语言处理任务，比如机器翻译和文本生成。

因为只依赖之前的输出，因此需要将训练样本的后面部分进行遮挡，这个过程称为掩码。

掩码具体的实现方式在注意力计算部分，在softmax时将掩码部分权重置0。

## 4.3 Transformer模型架构 【解码器】

- 1) 对输入进行掩码自注意力计算，并进行残差和层标准化计算，得到 $q$ 。
- 2) 将 $q$ 和Encoder部分的输出向量中的 $k$ 和 $v$ 一起进行自注意力计算，并进行残差和层标准化，得到 $v$ 。





softmax输出

$$\hat{y}=(\hat{y}^1, \hat{y}^2, \hat{y}^3, ..., \hat{y}_{\text{vocabulary\_size}})$$

线性层

$$d_{\text{model}} * \text{vocabulary\_size}$$

解码器输出

$$z=(z_1, z_2, z_3, ..., z_{d_{\text{model}}})$$

## 4.4 Transformer训练过程 【参数】

训练数据和批次	主要参数	硬件和训练时间	优化器
1) 标准的WMT2014英德数据集。	1) 模型维度: 512	1) 一台装有8个 NVIDIA P100 GPU 的机器。	1) Adam优化器: $\beta_1 = 0.9$ , $\beta_2 = 0.98$ 和 $\epsilon = 10^{-9}$ 。
2) 包含约450万个句子对。	2) 头数: 8	2) 一轮大约需要0.4秒。	2) warmup_steps个训练步骤中线性增加学习率, 之后按步骤数的倒数平方成比例地降低。 warmup_steps=4000
3) byte-pair 编码对句子进行编码。	3) 多头部分q、k v长度: $512/8 = 64$	3) 共10万步或12小时。	
4) 源-目标词汇表共包含约37000个标记。	4) dropout: 0.1		
5) 每个训练批次包含一组成约25000个源标记和25000个目标标记的句子对。	5) 标签平滑: $\epsilon_{ls} = 0.1$		

## 4.4 Transformer训练过程 【结果】

Transformer 在英语到德语和英语到法语 newstest2014 测试中取得了比以前最先进的模型更好的 BLEU 得分，但训练成本仅是之前的一小部分

Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [18]	23.75			
Deep-Att + PosUnk [39]		39.2		$1.0 \cdot 10^{20}$
GNMT + RL [38]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [9]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [32]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble [39]		40.4		$8.0 \cdot 10^{20}$
GNMT + RL Ensemble [38]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [9]	26.36	<b>41.29</b>	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (base model)	27.3	38.1	<b><math>3.3 \cdot 10^{18}</math></b>	
Transformer (big)	<b>28.4</b>	<b>41.8</b>	$2.3 \cdot 10^{19}$	

## 4.5 Transformer代码实操 【整体架构】

哈佛NLP团队复现了Transformer，地址：<https://github.com/harvardnlp/annotated-transformer>。

### 1) 原始输入部分：

```
nn.Sequential(Embeddings(d_model, src_vocab), c(position))
```

### 2) 编码器部分：

```
Encoder(EncoderLayer(d_model, c(attn), c(ff), dropout), N)
```

### 3) 目标输入部分：

```
nn.Sequential(Embeddings(d_model, tgt_vocab), c(position))
```

### 4) 解码器部分：

```
Decoder(DecoderLayer(d_model, c(attn), c(attn), c(ff), dropout), N)
```

### 5) 目标输出部分：

```
Generator(d_model, tgt_vocab)
```

```
def make_model(
    src_vocab, tgt_vocab, N=6, d_model=512, d_ff=2048, h=8, dropout=0.1
):
    "Helper: Construct a model from hyperparameters."
    c = copy.deepcopy
    attn = MultiHeadedAttention(h, d_model)
    ff = PositionwiseFeedForward(d_model, d_ff, dropout)
    position = PositionalEncoding(d_model, dropout)
    model = EncoderDecoder(
        Encoder(EncoderLayer(d_model, c(attn), c(ff), dropout), N),
        Decoder(DecoderLayer(d_model, c(attn), c(attn), c(ff), dropout), N),
        nn.Sequential(Embeddings(d_model, src_vocab), c(position)),
        nn.Sequential(Embeddings(d_model, tgt_vocab), c(position)),
        Generator(d_model, tgt_vocab),
    )

    # Initialize parameters with Glorot / fan_avg.
    for p in model.parameters():
        if p.dim() > 1:
            nn.init.xavier_uniform_(p)
    return model
```

## 4.5 Transformer代码实操 【训练概要】

EncoderDecoder: 负责实现编码器和解码器功能。

Generator: 负责实现最终输出部分功能。

```
class Generator(nn.Module):
    "Define standard linear + softmax generation step."

    def __init__(self, d_model, vocab):
        super(Generator, self).__init__()
        # See: https://pytorch.org/docs/stable/generated/torch.nn.Linear.html
        self.proj = nn.Linear(d_model, vocab)

    def forward(self, x):
        return log_softmax(self.proj(x), dim=-1)
```

```
class EncoderDecoder(nn.Module):
    """
    A standard Encoder-Decoder architecture. Base for this and many
    other models.
    """

    def __init__(self, encoder, decoder, src_embed, tgt_embed, generator):
        super(EncoderDecoder, self).__init__()
        self.encoder = encoder
        self.decoder = decoder
        self.src_embed = src_embed
        self.tgt_embed = tgt_embed
        self.generator = generator

    def forward(self, src, tgt, src_mask, tgt_mask):
        "Take in and process masked src and target sequences."
        return self.decode(self.encode(src, src_mask), src_mask, tgt, tgt_mask)

    def encode(self, src, src_mask):
        return self.encoder(self.src_embed(src), src_mask)

    def decode(self, memory, src_mask, tgt, tgt_mask):
        return self.decoder(self.tgt_embed(tgt), memory, src_mask, tgt_mask)
```

```
class Embeddings(nn.Module):
    def __init__(self, d_model, vocab):
        super(Embeddings, self).__init__()
        self.lut = nn.Embedding(vocab, d_model)
        self.d_model = d_model

    def forward(self, x):
        return self.lut(x) * math.sqrt(self.d_model)
```



负责Embedding层功能

```
class PositionalEncoding(nn.Module):
    "Implement the PE function."

    def __init__(self, d_model, dropout, max_len=5000):
        super(PositionalEncoding, self).__init__()
        self.dropout = nn.Dropout(p=dropout)

        # Compute the positional encodings once in log space.
        pe = torch.zeros(max_len, d_model)
        position = torch.arange(0, max_len).unsqueeze(1)
        div_term = torch.exp(
            torch.arange(0, d_model, 2) * -(math.log(10000.0) / d_model)
        )
        pe[:, 0::2] = torch.sin(position * div_term)
        pe[:, 1::2] = torch.cos(position * div_term)
        pe = pe.unsqueeze(0)
        self.register_buffer("pe", pe)

    def forward(self, x):
        x = x + self.pe[:, : x.size(1)].requires_grad_(False)
        return self.dropout(x)
```



负责位置编码功能

## 4.5 Transformer代码实操

### 【编码器部分一】

```
class Encoder(nn.Module):  
    "Core encoder is a stack of N layers"
```

负责编码器功能

```
    def __init__(self, layer, N):  
        super(Encoder, self).__init__()  
        self.layers = clones(layer, N)  
        self.norm = LayerNorm(layer.size)
```

```
    def forward(self, x, mask):  
        "Pass the input (and mask) through each layer in turn."  
        for layer in self.layers:  
            x = layer(x, mask)  
        return self.norm(x)
```

```
class EncoderLayer(nn.Module):  
    "Encoder is made up of self-attn and feed forward (defined below)"
```

负责编码器每一层功能

```
    def __init__(self, size, self_attn, feed_forward, dropout):  
        super(EncoderLayer, self).__init__()  
        self.self_attn = self_attn  
        self.feed_forward = feed_forward  
        self.sublayer = clones(SublayerConnection(size, dropout), 2)  
        self.size = size
```

```
    def forward(self, x, mask):  
        "Follow Figure 1 (left) for connections."  
        x = self.sublayer[0](x, lambda x: self.self_attn(x, x, x, mask))  
        return self.sublayer[1](x, self.feed_forward)
```

```
class SublayerConnection(nn.Module):  
    """
```

A residual connection followed by a layer norm.  
Note for code simplicity the norm is first as opposed to last.

负责层的连接：残差和层标准化功能

```
    def __init__(self, size, dropout):  
        super(SublayerConnection, self).__init__()  
        self.norm = LayerNorm(size)  
        self.dropout = nn.Dropout(dropout)
```

```
    def forward(self, x, sublayer):  
        "Apply residual connection to any sublayer with the same size."  
        return x + self.dropout(sublayer(self.norm(x)))
```

```
class PositionwiseFeedForward(nn.Module):  
    "Implements FFN equation."
```

负责基于位置的全连接功能

```
    def __init__(self, d_model, d_ff, dropout=0.1):  
        super(PositionwiseFeedForward, self).__init__()  
        self.w_1 = nn.Linear(d_model, d_ff)  
        self.w_2 = nn.Linear(d_ff, d_model)  
        self.dropout = nn.Dropout(dropout)
```

```
    def forward(self, x):  
        return self.w_2(self.dropout(self.w_1(x).relu()))
```

```
class MultiHeadedAttention(nn.Module):  
    def __init__(self, h, d_model, dropout=0.1):  
        "Take in model size and number of heads."  
        super(MultiHeadedAttention, self).__init__()  
        assert d_model % h == 0
```

负责多头注意力功能

```
        # We assume d_v always equals d_k  
        self.d_k = d_model // h  
        self.h = h  
        self.linears = clones(nn.Linear(d_model, d_model), 4)  
        self.attn = None  
        self.dropout = nn.Dropout(p=dropout)  
  
    def forward(self, query, key, value, mask=None):  
        "Implements Figure 2"  
        if mask is not None:  
            # Same mask applied to all h heads.  
            mask = mask.unsqueeze(1)  
            nbatches = query.size(0)  
  
        # 1) Do all the linear projections in batch from d_model => h x d_k  
        query, key, value = [  
            lin(x).view(nbatches, -1, self.h, self.d_k).transpose(1, 2)  
            for lin, x in zip(self.linears, (query, key, value))  
        ]  
  
        # 2) Apply attention on all the projected vectors in batch.  
        x, self.attn = attention(  
            query, key, value, mask=mask, dropout=self.dropout  
        )  
  
        # 3) "Concat" using a view and apply a final linear.  
        x = (  
            x.transpose(1, 2)  
            .contiguous()  
            .view(nbatches, -1, self.h * self.d_k)  
        )  
        del query  
        del key  
        del value  
        return self.linears[-1](x)
```

接上页

```
class LayerNorm(nn.Module):  
    "Construct a layernorm module (See citation for details)."  
  
    def __init__(self, features, eps=1e-6):  
        super(LayerNorm, self).__init__()  
        self.a_2 = nn.Parameter(torch.ones(features))  
        self.b_2 = nn.Parameter(torch.zeros(features))  
        self.eps = eps  
  
    def forward(self, x):  
        mean = x.mean(-1, keepdim=True)  
        std = x.std(-1, keepdim=True)  
        return self.a_2 * (x - mean) / (std + self.eps) + self.b_2
```

负责层标准化功能

```
def attention(query, key, value, mask=None, dropout=None):  
    "Compute 'Scaled Dot Product Attention'"  
    d_k = query.size(-1)  
    # Here is the attention function  
    scores = torch.matmul(query, key.transpose(-2, -1)) / math.sqrt(d_k)  
    if mask is not None:  
        scores = scores.masked_fill(mask == 0, -1e9)  
    p_attn = scores.softmax(dim=-1)  
    if dropout is not None:  
        p_attn = dropout(p_attn)  
    return torch.matmul(p_attn, value), p_attn
```

负责自注意力计算功能



```
class Decoder(nn.Module):  
    "Generic N layer decoder with masking."
```

```
    def __init__(self, layer, N):  
        super(Decoder, self).__init__()  
        self.layers = clones(layer, N)  
        self.norm = LayerNorm(layer.size)
```

```
    def forward(self, x, memory, src_mask, tgt_mask):  
        for layer in self.layers:  
            x = layer(x, memory, src_mask, tgt_mask)  
        return self.norm(x)
```



负责解码器部分功能

```
class DecoderLayer(nn.Module):  
    "Decoder is made of self-attn, src-attn, and feed forward (defined below)"
```

```
    def __init__(self, size, self_attn, src_attn, feed_forward, dropout):  
        super(DecoderLayer, self).__init__()  
        self.size = size  
        self.self_attn = self_attn  
        self.src_attn = src_attn  
        self.feed_forward = feed_forward  
        self.sublayer = clones(SublayerConnection(size, dropout), 3)
```

```
    def forward(self, x, memory, src_mask, tgt_mask):  
        "Follow Figure 1 (right) for connections."  
        m = memory  
        x = self.sublayer[0](x, lambda x: self.self_attn(x, x, x, tgt_mask))  
        x = self.sublayer[1](x, lambda x: self.src_attn(x, m, m, src_mask))  
        return self.sublayer[2](x, self.feed_forward)
```



负责解码器每一层功能

感谢您的观看

