



# Transformer

赵洲

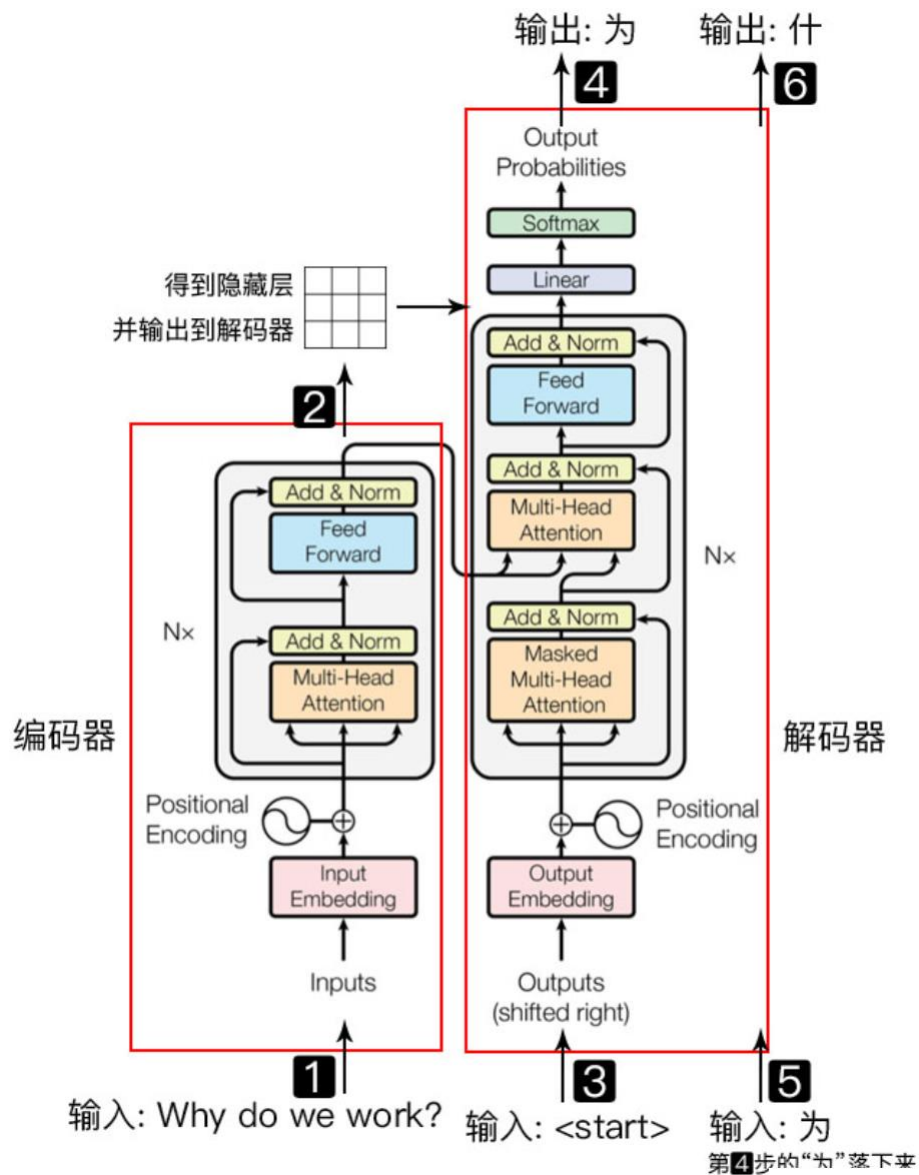
浙江大学计算机学院

2022年10月18日

# 模型结构

- 目前大部分比较热门的神经序列转换模型都有Encoder-Decoder结构。Encoder将输入序列 映射到一个连续表示序列
- 对于编码得到的，Decoder每次解码生成一个符号，直到生成完整的输出序列
- 对于每一步解码，模型都是自回归的，即在生成下一个符号时将先前生成的符号作为附加输入

# 模型结构图



# Attention机制

- 最简单的是QKV模型，假设输入为 $q$ , Memory中以 $(k, v)$ 形式储存需要的上下文。
- 如， $k$ 是question,  $v$ 是answer,  $q$ 是新来的question, 看看历史中memory中的哪个 $k$ 更相似，根据相似 $k$ 对应的 $v$ ，就是当前的答案。
  - 在memory中找相似 (score function):  $e_i = a(q, k_i)$
  - 归一化 (alignment function):  $\alpha = \text{softmax}(e_i)$
  - 读取内容 (context vector function):  $c = \sum_i \alpha_i v_i$

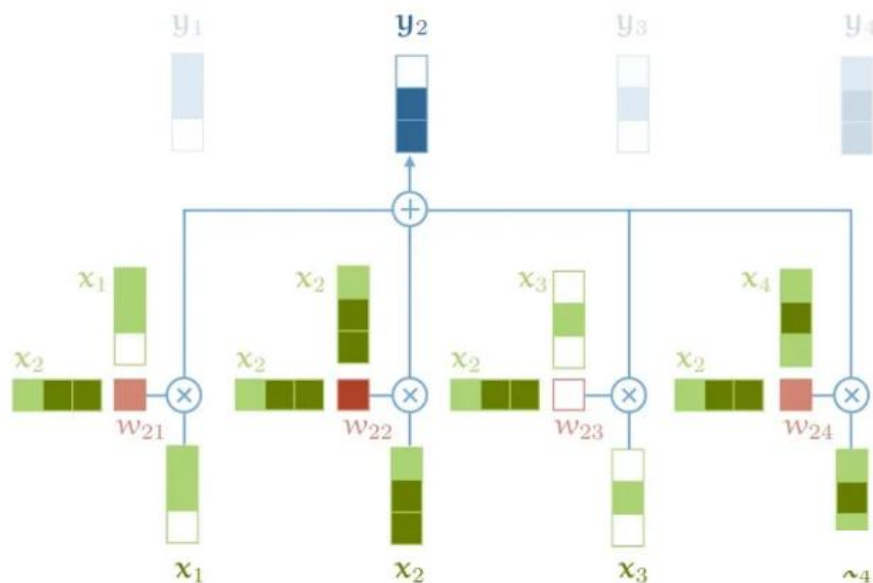
# Self-Attention机制

- Self-attention本质上是一种特殊的attention。
- 权重 $w$ 并不是一个需要神经网络学习的参与，来源于 $x_i$ 和 $x_j$ 的计算结果。

$$w'_{ij} = x_i^T x_j.$$

- 由于点积输出的取值范围是负无穷和正无穷之间，通过softmax映射到 $[0,1]$

$$w_{ij} = \frac{\exp w'_{ij}}{\sum_j \exp w'_{ij}}$$



# QKV聚焦机制

- Attention函数可以将Query和一组Key-Value对映射到输出，其中Query、Key、Value和输出都是向量。输出是值的加权和，其中分配给每个Value的权重由Query与相应Key的兼容函数计算
- 称这种特殊的Attention机制为"Scaled Dot-Product Attention"。输入包含维度为的Query和Key，以及维度为的Value

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

- 首先分别计算Query与各个Key的点积，然后将每个点积除以，最后使用Softmax函数来获得Key的权重。

# Self-Attention实现

```
1 import torch
2 import torch.nn.functional as F
3
4 # assume we have some tensor x with size (b, t, k)
5 x = ...
6
7 raw_weights = torch.bmm(x, x.transpose(1, 2))
8 # - torch.bmm is a batched matrix multiplication. It
9 #   applies matrix multiplication over batches of
10 #   matrices.
```

```
1 weights = F.softmax(raw_weights, dim=2)
```

```
1 y = torch.bmm(weights, x)
```

# Self-Attention的Trick 1

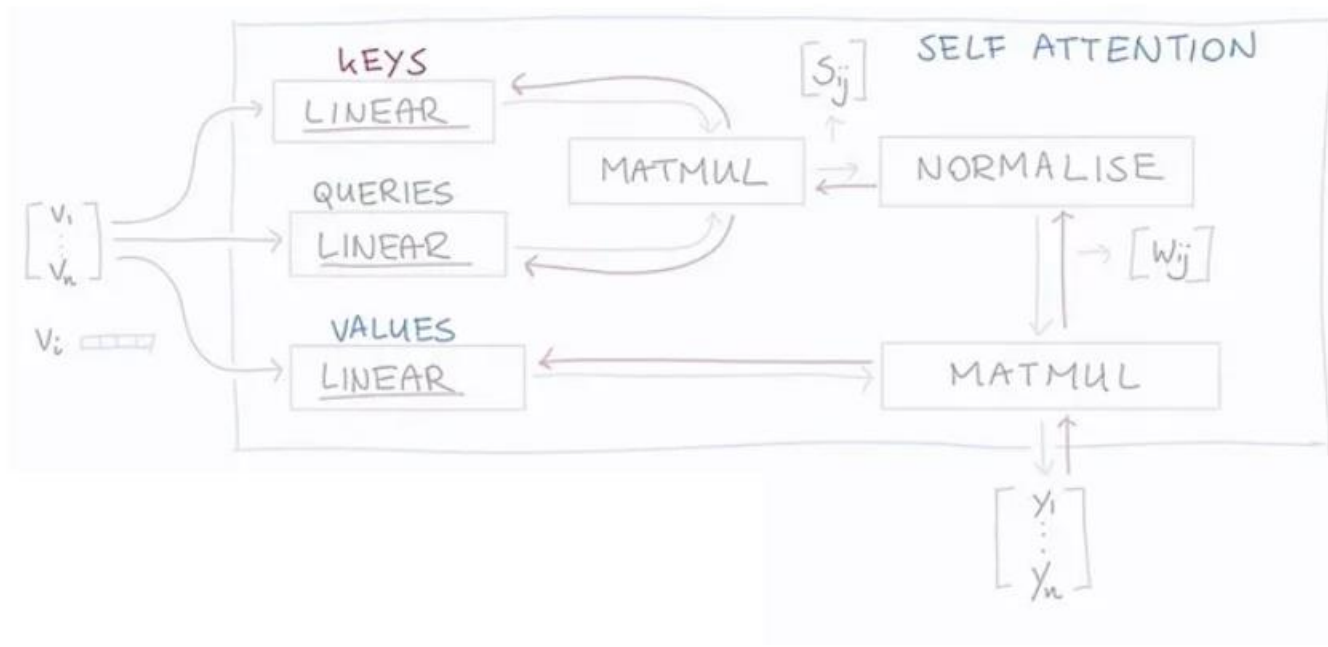
- 针对SA， 加入可学习参数

$$\mathbf{q}_i = \mathbf{W}_q \mathbf{x}_i \quad \mathbf{k}_i = \mathbf{W}_k \mathbf{x}_i \quad \mathbf{v}_i = \mathbf{W}_v \mathbf{x}_i$$

$$w'_{ij} = \mathbf{q}_i^T \mathbf{k}_j$$

$$w_{ij} = \text{softmax}(w'_{ij})$$

$$\mathbf{y}_i = \sum_j w_{ij} \mathbf{v}_j .$$





# Self-Attention的Trick 2

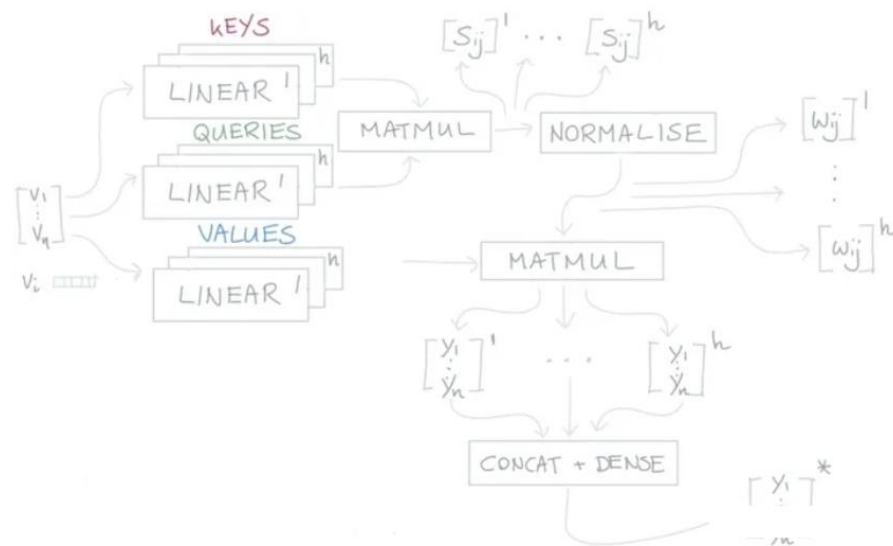
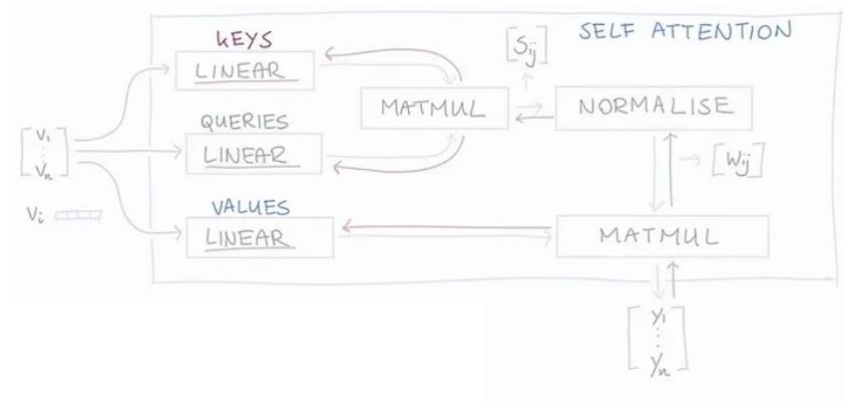
- Softmax函数对非常大的输入很敏感，会使得梯度传播出现问题 (kill the gradient)，导致学习的速度下降(slow down learning)，甚至导致学习的停止。
- 用 $\sqrt{k}$ 来对输入的向量做缩放，防止进入到softmax的函数增长过大

$$w'_{ij} = \frac{\mathbf{q}_i^T \mathbf{k}_j}{\sqrt{k}}$$

$$\begin{aligned}\frac{\partial a_j}{\partial z_l} &= \frac{\partial}{\partial z_l} \left( \frac{e^{z_j}}{\sum_k e^{z_k}} \right) \\ &= \frac{(e^{z_j})' \cdot \sum_k e^{z_k} - e^{z_j} \cdot e^{z_j}}{\left( \sum_k e^{z_k} \right)^2} \\ &= \frac{e^{z_j}}{\sum_k e^{z_k}} - \frac{e^{z_j}}{\sum_k e^{z_k}} \cdot \frac{e^{z_j}}{\sum_k e^{z_k}} = a_j(1 - a_j)\end{aligned}$$

# Self-Attention的Trick 3

- 对于输入 $x$ 每一个attention head都会生成一个向量 $y_i$ 。把这些向量进行concat操作，并且把concat的结果传递给一个全连接层，使得向量的维度重新回到 $k$ 。

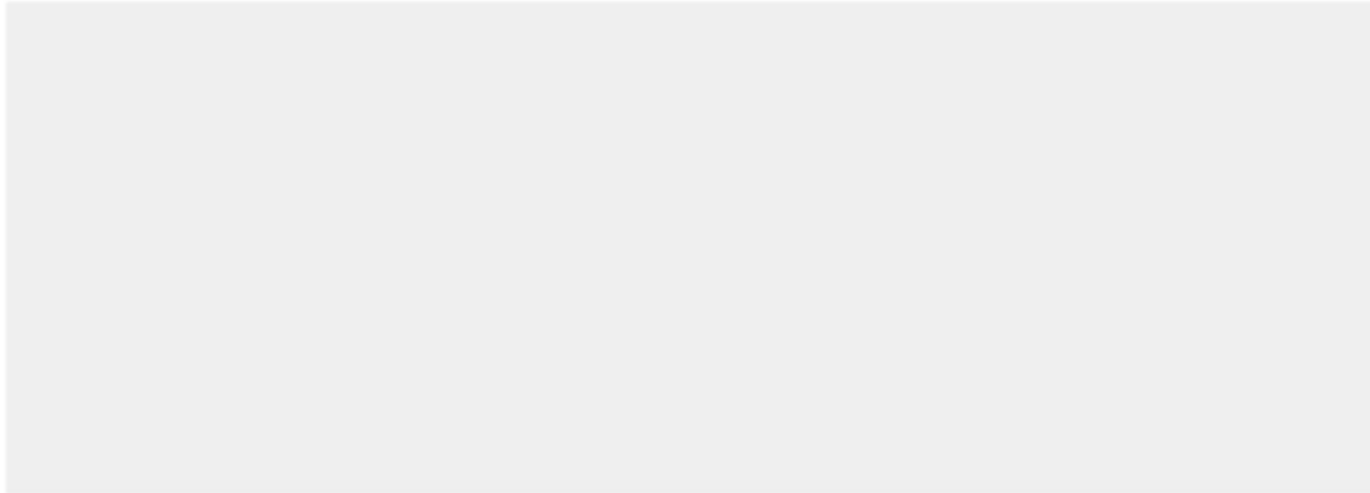


# Self-Attention实现例子

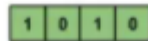
1. 准备输入
2. 初始化参数
3. 获取key, query和value
4. 给input1计算attention score
5. 给value乘上score
6. 给value加权求和获取output1
7. 重复步骤4-6, 获取output2, output3

# 准备输入

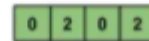
Self-attention



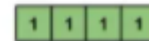
input #1



input #2



input #3



- 1 Input 1: [1, 0, 1, 0]
- 2 Input 2: [0, 2, 0, 2]
- 3 Input 3: [1, 1, 1, 1]

# 初始化参数

key的参数:

```
1 [[0, 0, 1],  
2  [1, 1, 0],  
3  [0, 1, 0],  
4  [1, 1, 0]]
```

query的参数:

```
1 [[1, 0, 1],  
2  [1, 0, 0],  
3  [0, 0, 1],  
4  [0, 1, 1]]
```

value的参数:

```
1 [[0, 2, 0],  
2  [0, 3, 0],  
3  [1, 0, 3],  
4  [1, 1, 0]]
```

# 获取key, query和value

对于input1的key的表示为:

```
1           [0, 0, 1]
2  [1, 0, 1, 0] × [1, 1, 0] = [0, 1, 1]
3           [0, 1, 0]
4           [1, 1, 0]
```

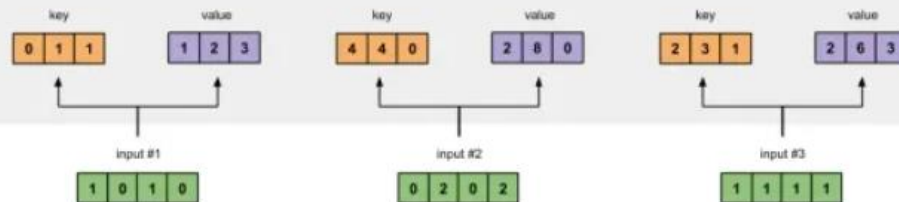
使用相同的参数获取input2的key的表示:

```
1           [0, 0, 1]
2  [0, 2, 0, 2] × [1, 1, 0] = [4, 4, 0]
3           [0, 1, 0]
4           [1, 1, 0]
```

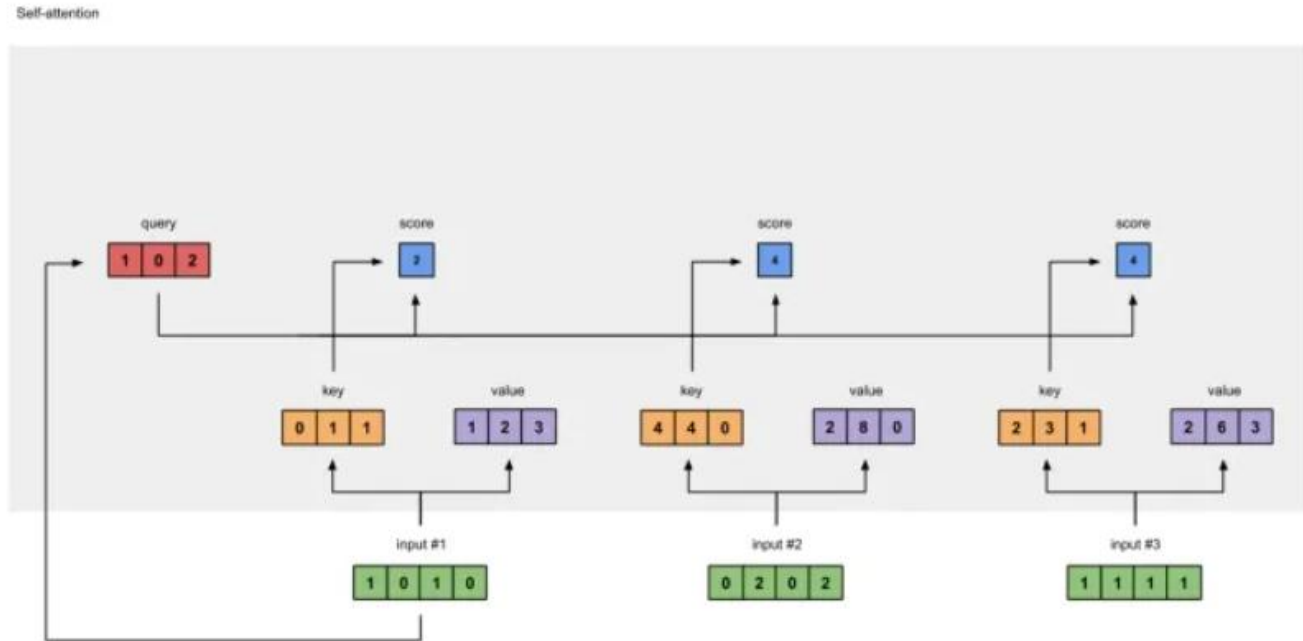
使用参数获取input3的key的表示:

```
1           [0, 0, 1]
2  [1, 1, 1, 1] × [1, 1, 0] = [2, 3, 1]
3           [0, 1, 0]
4           [1, 1, 0]
```

Self-attention



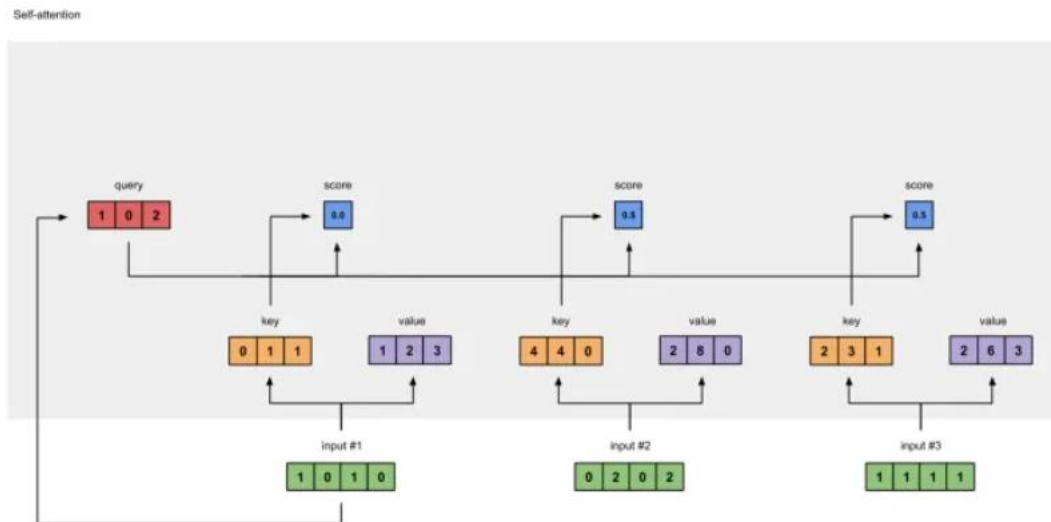
# 给input1计算attention score



给attention score应用softmax。

```
1 softmax([2, 4, 4]) = [0.0, 0.5, 0.5]
```

# 给value乘上score

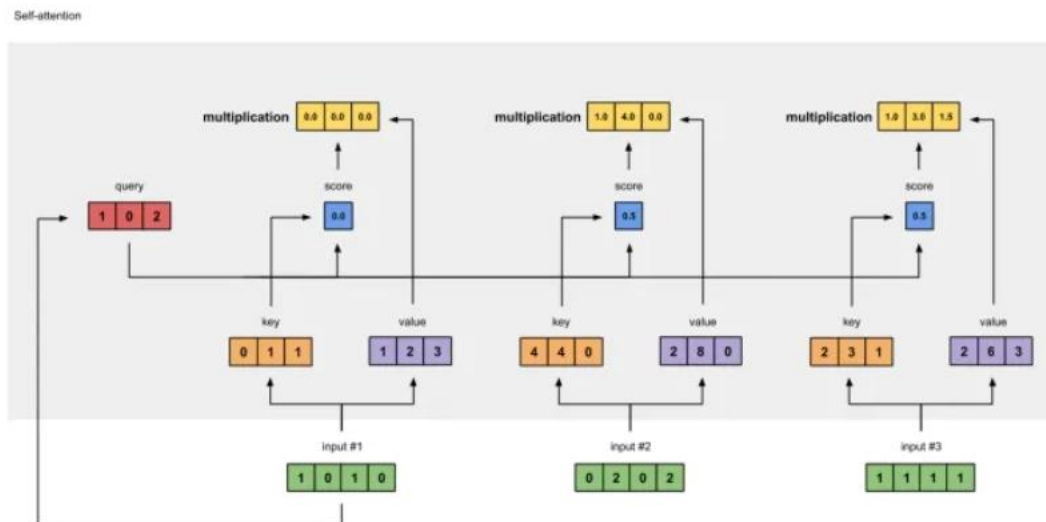


使用经过softmax后的attention score乘以它对应的value值（紫色），这样我们就得到了3个weighted values（黄色）。

- 1:  $0.0 * [1, 2, 3] = [0.0, 0.0, 0.0]$
- 2:  $0.5 * [2, 8, 0] = [1.0, 4.0, 0.0]$
- 3:  $0.5 * [2, 6, 3] = [1.0, 3.0, 1.5]$



# 给value加权求和获取output1



把所有的weighted values（黄色）进行element-wise的相加。

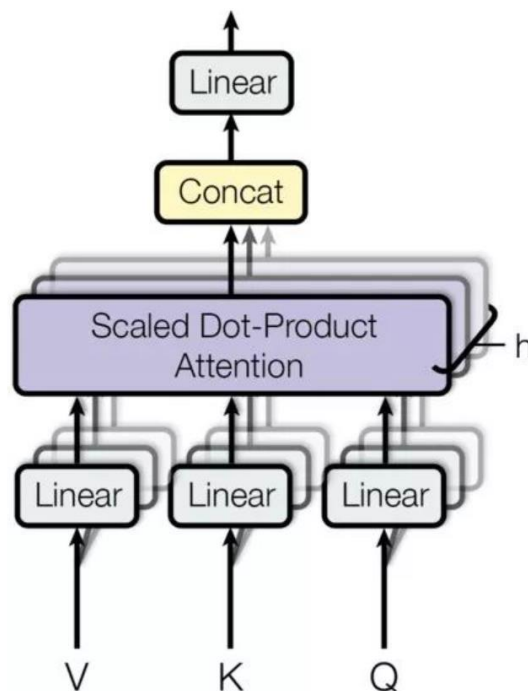
```
1  [0.0, 0.0, 0.0]
2  + [1.0, 4.0, 0.0]
3  + [1.0, 3.0, 1.5]
4  -----
5  = [2.0, 7.0, 1.5]
```

# 多头Attention

- “多头”机制能让模型考虑到不同位置的Attention，另外“多头”Attention可以在不同的子空间表示不一样的关联关系，使用单个Head的Attention一般达不到这种效果

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

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



# Transformer模型by Pytorch

- Tokenize
- Input Embedding
- Positional Encoder
- Transformer Block
- Encoder
- Decoder
- Transformer

# Tokenize

it not cool that ping pong is not included in rio 2016

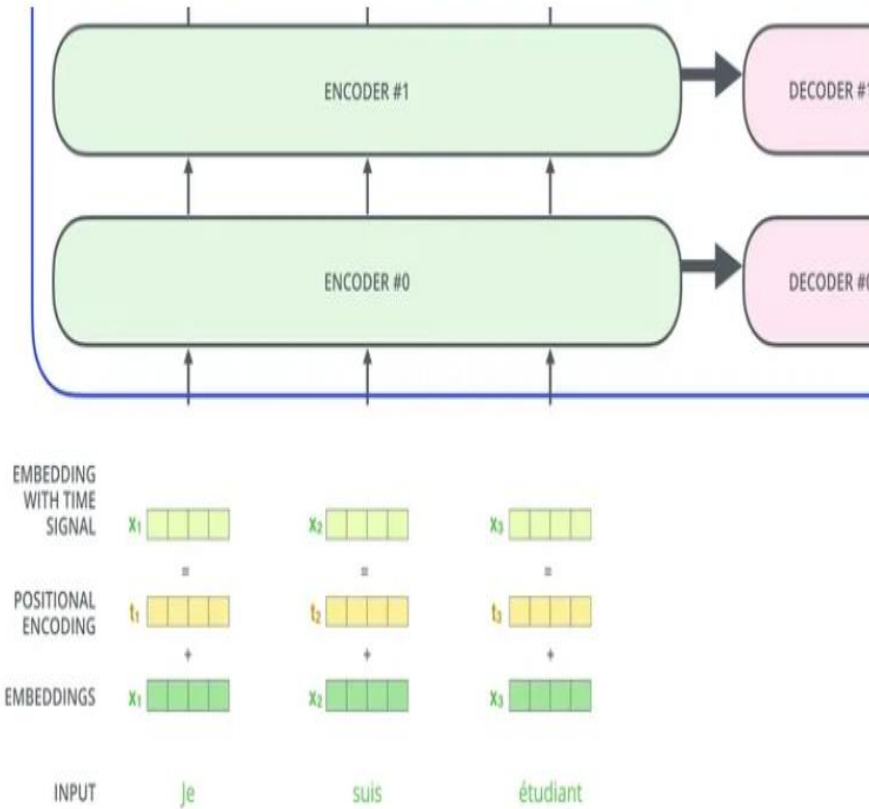


**Tokenization**



```
1 class Tokenize(object):
2
3     def __init__(self, lang):
4         self.nlp = importlib.import_module(lang).load()
5
6     def tokenizer(self, sentence):
7         sentence = re.sub(
8             r"[\*\"'\"'\n\\...\\+\\-\\/\\=\\(\\)‘•:;\\[\\]\\\\|'\\\\!;]", " ", str(sentence))
9         sentence = re.sub(r"[ ]+", " ", sentence)
10        sentence = re.sub(r"!+", "!", sentence)
11        sentence = re.sub(r"\\,", ",", sentence)
12        sentence = re.sub(r"\\?+", "?", sentence)
13        sentence = sentence.lower()
14        return [tok.text for tok in self.nlp.tokenizer(sentence) if tok.text !=
```

# Input Embedding

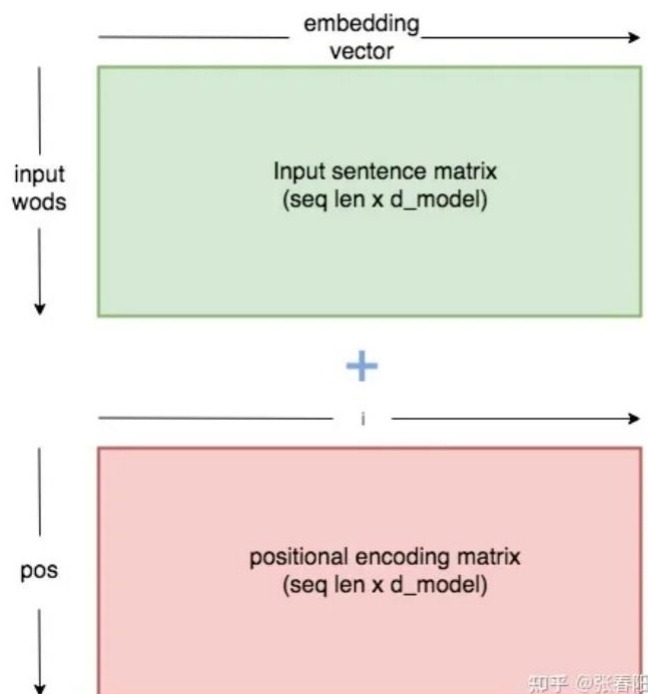


```
1 class Embedding(nn.Module):
2
3     def __init__(self, vocab_size, d_model):
4         super().__init__()
5         self.d_model = d_model
6         self.embed = nn.Embedding(vocab_size, d_model)
7
8     def forward(self, x):
9         return self.embed(x)
```

# Positional Encoder

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

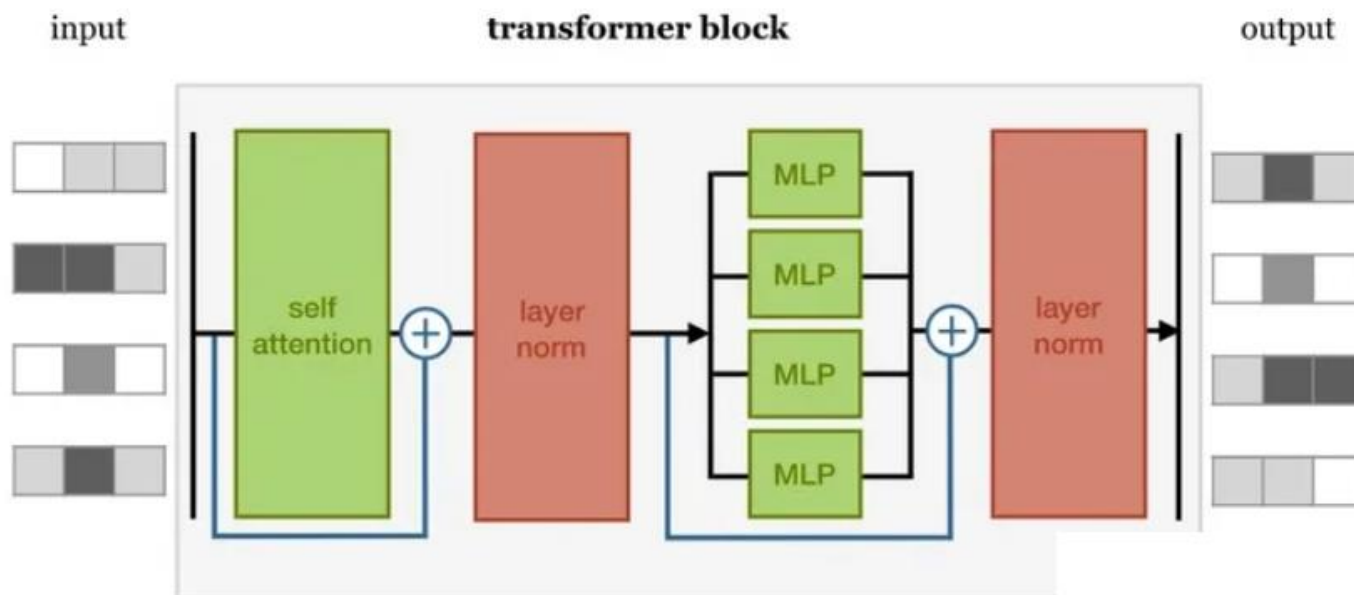
$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{model}})$$



```
1 class PositionalEncoder(nn.Module):
2
3     def __init__(self, d_model, max_seq_len = 80):
4         super().__init__()
5         self.d_model = d_model
6
7         # 根据pos和i创建一个常量pe矩阵
8         pe = torch.zeros(max_seq_len, d_model)
9         for pos in range(max_seq_len):
10             for i in range(0, d_model, 2):
11                 pe[pos, i] = \
12                     math.sin(pos / (10000 ** ((2 * i)/d_model)))
13                 pe[pos, i + 1] = \
14                     math.cos(pos / (10000 ** ((2 * (i + 1))/d_model)))
15
16         pe = pe.unsqueeze(0)
17         self.register_buffer('pe', pe)
18
19     def forward(self, x):
20         # 让 embeddings vector 相对大一些
21         x = x * math.sqrt(self.d_model)
22         # 增加位置常量到 embedding 中
23         seq_len = x.size(1)
24         x = x + Variable(self.pe[:, :seq_len], \
25                           requires_grad=False).cuda()
26
27         return x
```

# Transformer Block

- Transformer Block 主要是有以下4个部分构成的:
  - ◆ self-attention layer
  - ◆ normalization layer
  - ◆ feed forward layer
  - ◆ another normalization layer



# self-attention layer

```
1  def attention(q, k, v, d_k, mask=None, dropout=None):
2
3      scores = torch.matmul(q, k.transpose(-2, -1)) / math.sqrt(d_k)
4
5      # mask掉那些为了padding长度增加的token, 让其通过softmax计算后为0
6      if mask is not None:
7          mask = mask.unsqueeze(1)
8          scores = scores.masked_fill(mask == 0, -1e9)
9
10     scores = F.softmax(scores, dim=-1)
11
12     if dropout is not None:
13         scores = dropout(scores)
14
15     output = torch.matmul(scores, v)
16     return output
```



# MultiHead Attention

- 多头的注意力机制，用来识别数据之间的不同联系

```
1 class MultiHeadAttention(nn.Module):
2
3     def __init__(self, heads, d_model, dropout = 0.1):
4         super().__init__()
5
6         self.d_model = d_model
7         self.d_k = d_model // heads
8         self.h = heads
9
10        self.q_linear = nn.Linear(d_model, d_model)
11        self.v_linear = nn.Linear(d_model, d_model)
12        self.k_linear = nn.Linear(d_model, d_model)
13
14        self.dropout = nn.Dropout(dropout)
15        self.out = nn.Linear(d_model, d_model)
16
17    def forward(self, q, k, v, mask=None):
18
19        bs = q.size(0)
20
21        # perform linear operation and split into N heads
22        k = self.k_linear(k).view(bs, -1, self.h, self.d_k)
23        q = self.q_linear(q).view(bs, -1, self.h, self.d_k)
24        v = self.v_linear(v).view(bs, -1, self.h, self.d_k)
25
26        # transpose to get dimensions bs * N * sl * d_model
27        k = k.transpose(1,2)
28        q = q.transpose(1,2)
29        v = v.transpose(1,2)
30
31        # calculate attention using function we will define next
32        scores = attention(q, k, v, self.d_k, mask, self.dropout)
33        # concatenate heads and put through final linear layer
34        concat = scores.transpose(1,2).contiguous()\
35            .view(bs, -1, self.d_model)
36        output = self.out(concat)
37
38        return output
```

# Layer Norm

## ■ 使用 Layer Norm 来使得梯度更加的平稳

### Algorithm 1 Batch Normalization (Every Iteration)

#### begin Forward Propagation:

**Input:**  $X \in R^{B \times d}$

**Output:**  $Y \in R^{B \times d}$

$$\mu_B = \frac{1}{B} \sum_{i=1}^B \mathbf{x}_i \quad // \text{ Get mini-batch mean}$$

$$\sigma_B^2 = \frac{1}{B} \sum_{i=1}^B (\mathbf{x}_i - \mu_B)^2 \quad // \text{ Get mini-batch variance}$$

$$\tilde{X} = \frac{X - \mu_B}{\sigma_B} \quad // \text{ Normalize}$$

$$Y = \gamma \odot \tilde{X} + \beta \quad // \text{ Scale and shift}$$

$$\mu = \alpha \mu + (1 - \alpha) \mu_B \quad // \text{ Update running mean}$$

$$\sigma^2 = \alpha \sigma^2 + (1 - \alpha) \sigma_B^2 \quad // \text{ Update running variance}$$

#### begin Backward Propagation:

**Input:**  $\frac{\partial \mathcal{L}}{\partial Y} \in R^{B \times d}$

**Output:**  $\frac{\partial \mathcal{L}}{\partial X} \in R^{B \times d}$

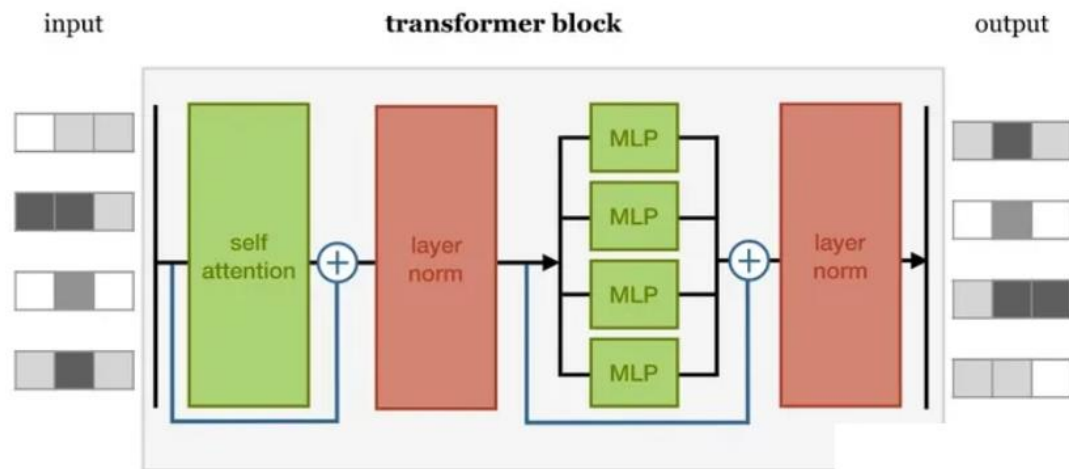
$$\frac{\partial \mathcal{L}}{\partial X} \text{ based on Eq. 3} \quad // \text{ Gradient of } X$$

$$\text{Inference: } Y = \gamma \odot \frac{X - \mu}{\sigma} + \beta$$

```
1 class NormLayer(nn.Module):
2
3     def __init__(self, d_model, eps = 1e-6):
4         super().__init__()
5
6         self.size = d_model
7
8         # 使用两个可以学习的参数来进行 normalisation
9         self.alpha = nn.Parameter(torch.ones(self.size))
10        self.bias = nn.Parameter(torch.zeros(self.size))
11
12        self.eps = eps
13
14    def forward(self, x):
15        norm = self.alpha * (x - x.mean(dim=-1, keepdim=True)) \
16        / (x.std(dim=-1, keepdim=True) + self.eps) + self.bias
17        return norm
```

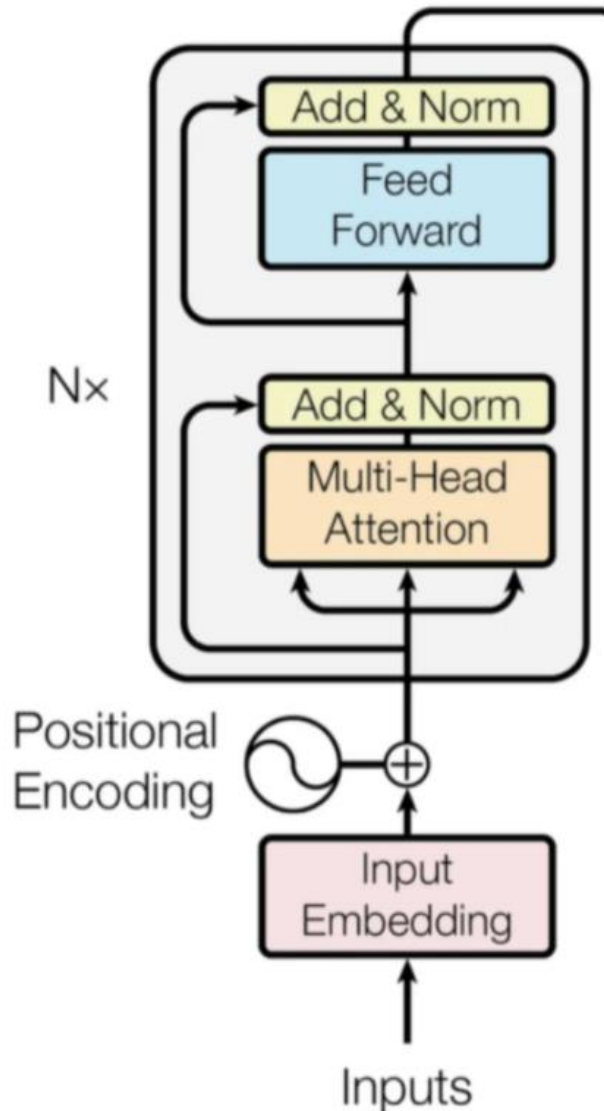
# Feed Forward Layer

- 选择4倍输入大小作为我们 feedforward 层的维度，这个值使用的越小就越节省内存，但是相应的表示性也会变弱



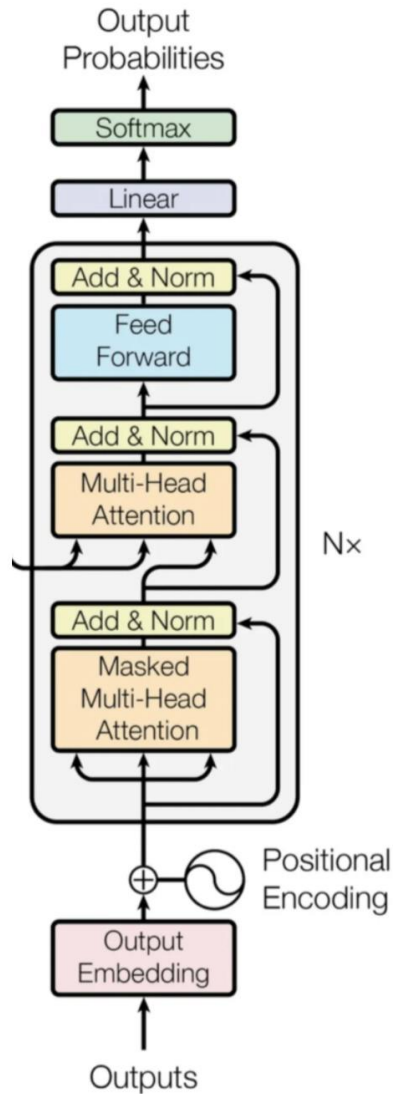
```
1 class TransformerBlock(nn.Module):
2     def __init__(self, k, heads):
3         super().__init__()
4
5         self.attention = SelfAttention(k, heads=heads)
6
7         self.norm1 = nn.LayerNorm(k)
8         self.norm2 = nn.LayerNorm(k)
9
10        self.ff = nn.Sequential(
11            nn.Linear(k, 4 * k),
12            nn.ReLU(),
13            nn.Linear(4 * k, k))
14
15    def forward(self, x):
16        attended = self.attention(x)
17        x = self.norm1(attended + x)
18
19        feedforward = self.ff(x)
20        return self.norm2(feedforward + x)
```

# Encoder



```
1 class EncoderLayer(nn.Module):
2
3     def __init__(self, d_model, heads, dropout=0.1):
4         super().__init__()
5         self.norm_1 = Norm(d_model)
6         self.norm_2 = Norm(d_model)
7         self.attn = MultiHeadAttention(heads, d_model, dropout=dropout)
8         self.ff = FeedForward(d_model, dropout=dropout)
9         self.dropout_1 = nn.Dropout(dropout)
10        self.dropout_2 = nn.Dropout(dropout)
11
12    def forward(self, x, mask):
13        x2 = self.norm_1(x)
14        x = x + self.dropout_1(self.attn(x2, x2, x2, mask))
15        x2 = self.norm_2(x)
16        x = x + self.dropout_2(self.ff(x2))
17        return x
18
19
20 class Encoder(nn.Module):
21
22    def __init__(self, vocab_size, d_model, N, heads, dropout):
23        super().__init__()
24        self.N = N
25        self.embed = Embedder(vocab_size, d_model)
26        self.pe = PositionalEncoder(d_model, dropout=dropout)
27        self.layers = get_clones(EncoderLayer(d_model, heads, dropout), N)
28        self.norm = Norm(d_model)
29
30    def forward(self, src, mask):
31        x = self.embed(src)
32        x = self.pe(x)
33        for i in range(self.N):
34            x = self.layers[i](x, mask)
35        return self.norm(x)
```

# Decoder



```

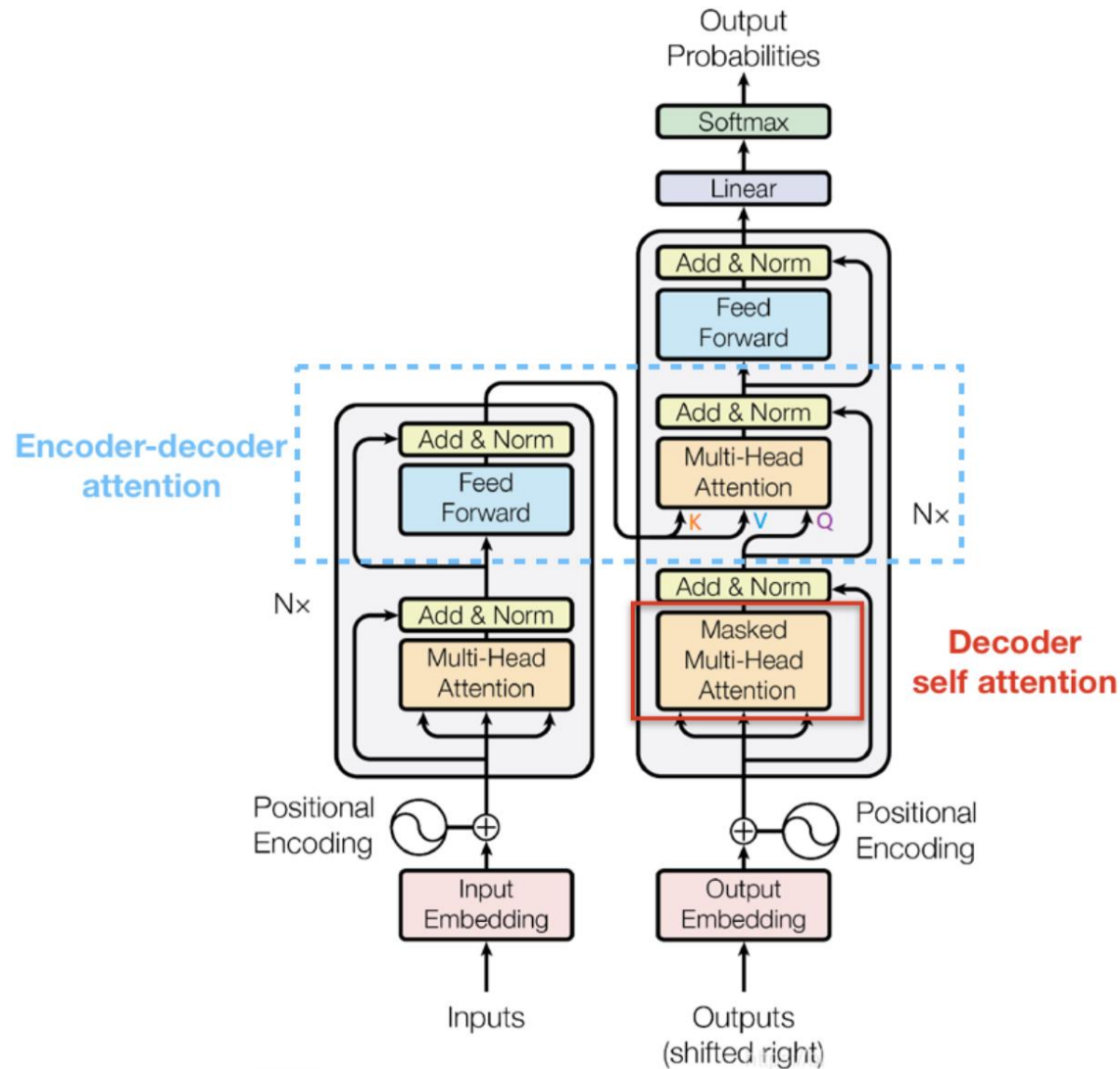
1 class DecoderLayer(nn.Module):
2
3     def __init__(self, d_model, heads, dropout=0.1):
4         super().__init__()
5         self.norm_1 = Norm(d_model)
6         self.norm_2 = Norm(d_model)
7         self.norm_3 = Norm(d_model)
8
9         self.dropout_1 = nn.Dropout(dropout)
10        self.dropout_2 = nn.Dropout(dropout)
11        self.dropout_3 = nn.Dropout(dropout)
12
13        self.attn_1 = MultiHeadAttention(heads, d_model, dropout=dropout)
14        self.attn_2 = MultiHeadAttention(heads, d_model, dropout=dropout)
15        self.ff = FeedForward(d_model, dropout=dropout)
16
17    def forward(self, x, e_outputs, src_mask, trg_mask):
18        x2 = self.norm_1(x)
19        x = x + self.dropout_1(self.attn_1(x2, x2, x2, trg_mask))
20        x2 = self.norm_2(x)
21        x = x + self.dropout_2(self.attn_2(x2, e_outputs, e_outputs, \
22        src_mask))
23        x2 = self.norm_3(x)
24        x = x + self.dropout_3(self.ff(x2))
25        return x
26
27
28 class Decoder(nn.Module):
29
30     def __init__(self, vocab_size, d_model, N, heads, dropout):
31         super().__init__()
32         self.N = N
33         self.embed = Embedder(vocab_size, d_model)
34         self.pe = PositionalEncoder(d_model, dropout=dropout)
35         self.layers = get_clones(DecoderLayer(d_model, heads, dropout), N)
36         self.norm = Norm(d_model)
37
38    def forward(self, trg, e_outputs, src_mask, trg_mask):
39        x = self.embed(trg)
40        x = self.pe(x)
41        for i in range(self.N):
42            x = self.layers[i](x, e_outputs, src_mask, trg_mask)
43        return self.norm(x)

```

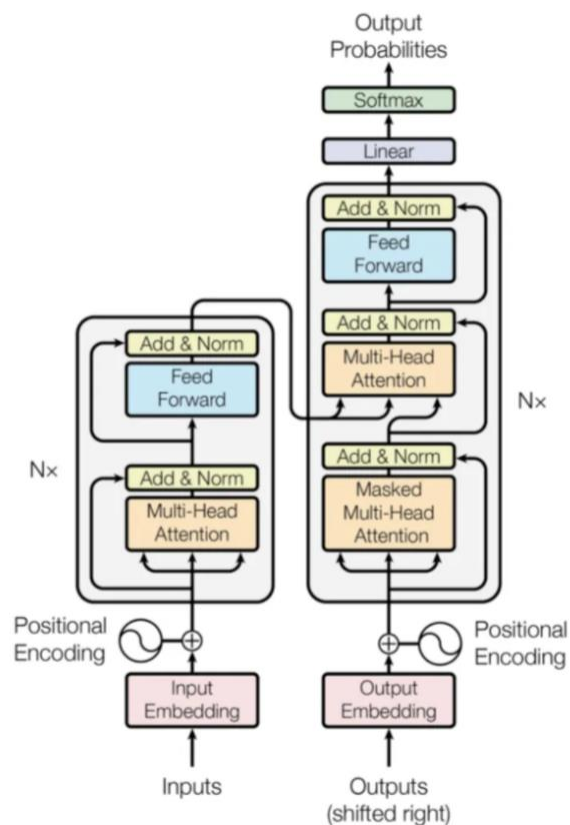
# Attention在Transformer中作用

- Encoder self-attention: Encoder 阶段捕获当前 word 和其他输入词的关联
- Masked-Decoder self-attention : Decoder 阶段捕获当前 word 与已经看到的解码词之间的关联，从矩阵上直观来看就是一个带有 mask 的三角矩阵
- Encoder-Decoder Attention: 就是将 Decoder 和 Encoder 输入建立联系，和之前那些普通 Attention 一样

# Attention in Decoder



# Transformer

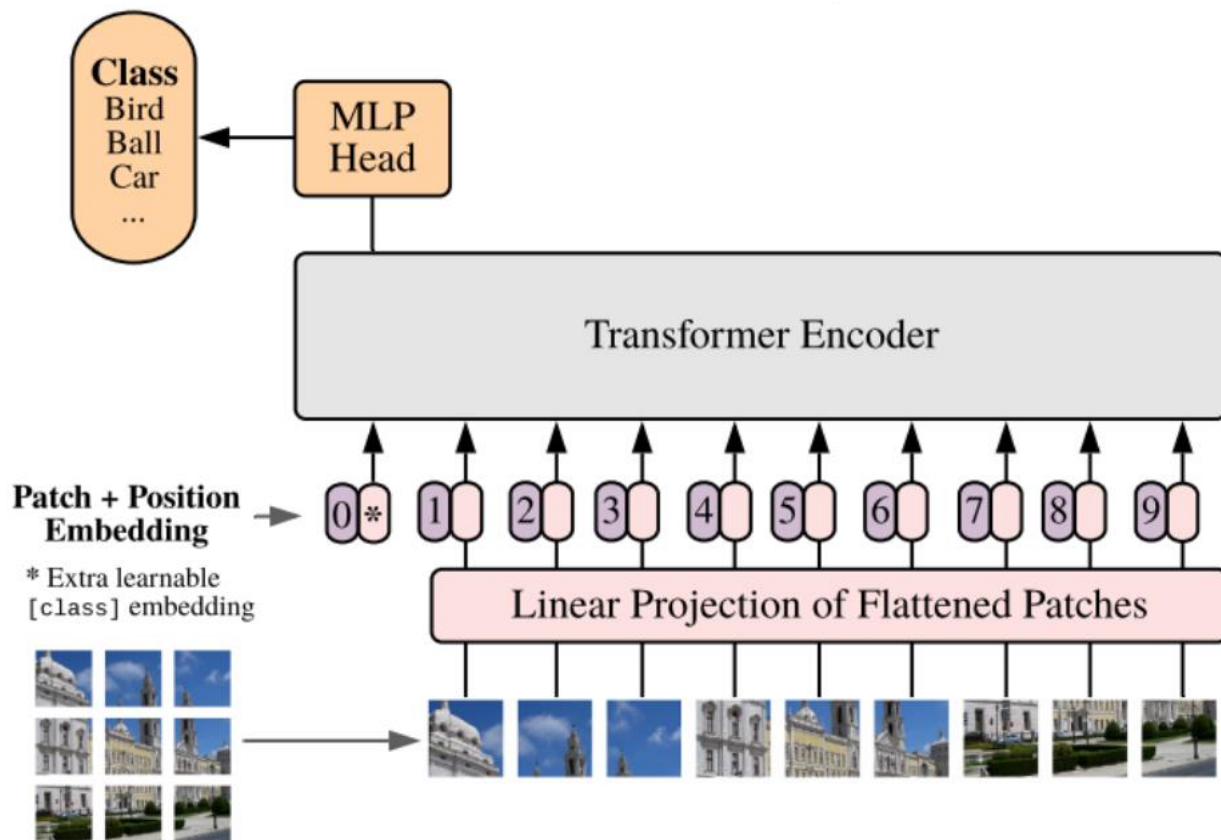


```
1 class Transformer(nn.Module):
2
3     def __init__(self, src_vocab, trg_vocab, d_model, N, heads, dropout):
4         super().__init__()
5         self.encoder = Encoder(src_vocab, d_model, N, heads, dropout)
6         self.decoder = Decoder(trg_vocab, d_model, N, heads, dropout)
7         self.out = nn.Linear(d_model, trg_vocab)
8
9     def forward(self, src, trg, src_mask, trg_mask):
10         e_outputs = self.encoder(src, src_mask)
11         d_output = self.decoder(trg, e_outputs, src_mask, trg_mask)
12         output = self.out(d_output)
13         return output
```

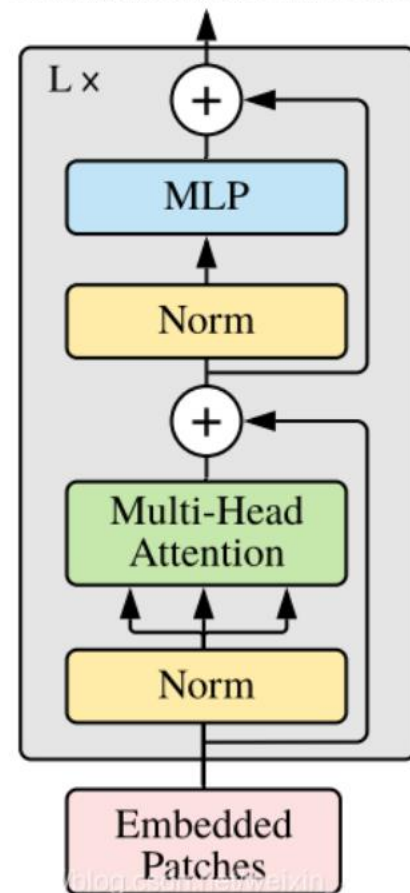


# Vision Transformer 整体结构

Vision Transformer (ViT)



Transformer Encoder



# 数据处理

- 原始输入的图片数据是 $H*W*C$ , 先对图片做分块, 再进行展平。  
假设每个块的长度为 $(P,P)$ , 那么分块的数目为:

$$N = H * W / (P * P)$$

- 每个图片块展平成一维向量, 每个向量大小为:

$$P * P * C$$

- 总的输入变换为:

$$N \times (P^2 \cdot C)$$

# Patch Embedding

- 对每个向量都做一个线性变换（即全连接层），压缩维度为D，称其为Patch Embedding
- 在代码里初始化一个全连接层，输出维度为dim, 然后将分块后的数据输入

```
1 self.patch_to_embedding = nn.Linear(patch_dim, dim)
2
3 # forward前向代码
4 x = rearrange(img, 'b c (h p1) (w p2) -> b (h w) (p1 p2 c)', p1=p, p2=p)
5 x = self.patch_to_embedding(x)
```

# Position Encoding

- 位置编码并没有使用传统的Transformer的cos/sin的那套编码方式

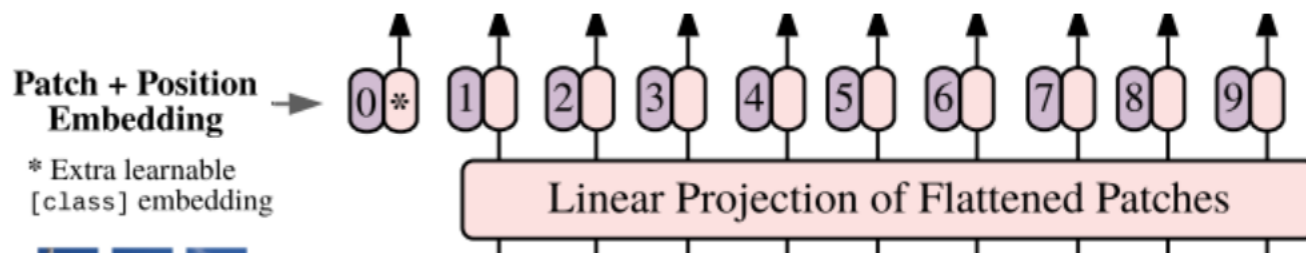
$$\begin{cases} PE(pos, 2i) = \sin(pos/10000^{2i/d_{\text{model}}}) \\ PE(pos, 2i + 1) = \cos(pos/10000^{2i/d_{\text{model}}}) \end{cases}$$

- 而是采用随机初始化，之后再训练出来

```
self.pos_embedding = nn.Parameter(torch.randn(1, num_patches + 1, dim))  
x += self.pos_embedding[:, :(n + 1)]
```

# Cls Token (分类)

- ViT只用了Encoder编码器结构，缺少解码过程，因此给了额外的一个用于分类的向量，与输入进行拼接



```
1 # 假设dim=128, 这里shape为(1, 1, 128)
2 self.cls_token = nn.Parameter(torch.randn(1, 1, dim))
3
4 # forward前向代码
5 # 假设batchsize=10, 这里shape为(10, 1, 128)
6 cls_tokens = repeat(self.cls_token, '() n d -> b n d', b=b)
7 # 跟前面的分块为x (10, 64, 128) 的进行concat
8 # 得到 (10, 65, 128) 向量
9 x = torch.cat((cls_tokens, x), dim=1)
```

# Transformer编码器

- Transformer编码器由多个交互层的多头自注意力和MLP块组成。
- 每个块之前应用LayerNorm, 残差连接在每个块之后应用。
- MLP包含两个呈现GELU非线性层。

$$\begin{aligned} \mathbf{z}_0 &= [\mathbf{x}_{\text{class}}; \mathbf{x}_p^1 \mathbf{E}; \mathbf{x}_p^2 \mathbf{E}; \cdots; \mathbf{x}_p^N \mathbf{E}] + \mathbf{E}_{\text{pos}}, & \mathbf{E} \in \mathbb{R}^{(P^2 \cdot C) \times D}, \mathbf{E}_{\text{pos}} \in \mathbb{R}^{(N+1) \times D} \\ \mathbf{z}'_\ell &= \text{MSA}(\text{LN}(\mathbf{z}_{\ell-1})) + \mathbf{z}_{\ell-1}, & \ell = 1 \dots L \\ \mathbf{z}_\ell &= \text{MLP}(\text{LN}(\mathbf{z}'_\ell)) + \mathbf{z}'_\ell, & \ell = 1 \dots L \\ \mathbf{y} &= \text{LN}(\mathbf{z}_L^0) \end{aligned}$$

# 分类

- 分类头加入LayerNorm和两层全连接层实现，采用GELU激活函数

```
self.mlp_head = nn.Sequential(  
    nn.LayerNorm(dim),  
    nn.Linear(dim, mlp_dim),  
    nn.GELU(),  
    nn.Dropout(dropout),  
    nn.Linear(mlp_dim, num_classes)  
)
```

- 最终分类我们只取第一个，也就是用于分类的token，输入到分类头里，得到分类结果

```
self.to_cls_token = nn.Identity()  
  
# forward前向部分  
  
x = self.transformer(x, mask)  
x = self.to_cls_token(x[:, 0])  
  
return self.mlp_head(x)
```

# 实验部分

| Model     | Layers | Hidden size $D$ | MLP size | Heads | Params |
|-----------|--------|-----------------|----------|-------|--------|
| ViT-Base  | 12     | 768             | 3072     | 12    | 86M    |
| ViT-Large | 24     | 1024            | 4096     | 16    | 307M   |
| ViT-Huge  | 32     | 1280            | 5120     | 16    | 632M   |

Table 1: Details of Vision Transformer model variants.

|                    | Ours-JFT<br>(ViT-H/14)  | Ours-JFT<br>(ViT-L/16)  | Ours-I21K<br>(ViT-L/16) | BiT-L<br>(ResNet152x4) | Noisy Student<br>(EfficientNet-L2) |
|--------------------|-------------------------|-------------------------|-------------------------|------------------------|------------------------------------|
| ImageNet           | <b>88.55</b> $\pm 0.04$ | 87.76 $\pm 0.03$        | 85.30 $\pm 0.02$        | 87.54 $\pm 0.02$       | 88.4/88.5*                         |
| ImageNet ReaL      | <b>90.72</b> $\pm 0.05$ | 90.54 $\pm 0.03$        | 88.62 $\pm 0.05$        | 90.54                  | 90.55                              |
| CIFAR-10           | <b>99.50</b> $\pm 0.06$ | 99.42 $\pm 0.03$        | 99.15 $\pm 0.03$        | 99.37 $\pm 0.06$       | —                                  |
| CIFAR-100          | <b>94.55</b> $\pm 0.04$ | 93.90 $\pm 0.05$        | 93.25 $\pm 0.05$        | 93.51 $\pm 0.08$       | —                                  |
| Oxford-IIIT Pets   | <b>97.56</b> $\pm 0.03$ | 97.32 $\pm 0.11$        | 94.67 $\pm 0.15$        | 96.62 $\pm 0.23$       | —                                  |
| Oxford Flowers-102 | 99.68 $\pm 0.02$        | <b>99.74</b> $\pm 0.00$ | 99.61 $\pm 0.02$        | 99.63 $\pm 0.03$       | —                                  |
| VTAB (19 tasks)    | <b>77.63</b> $\pm 0.23$ | 76.28 $\pm 0.46$        | 72.72 $\pm 0.21$        | 76.29 $\pm 1.70$       | —                                  |
| TPUv3-core-days    | 2.5k                    | 0.68k                   | 0.23k                   | 9.9k                   | 12.3k                              |