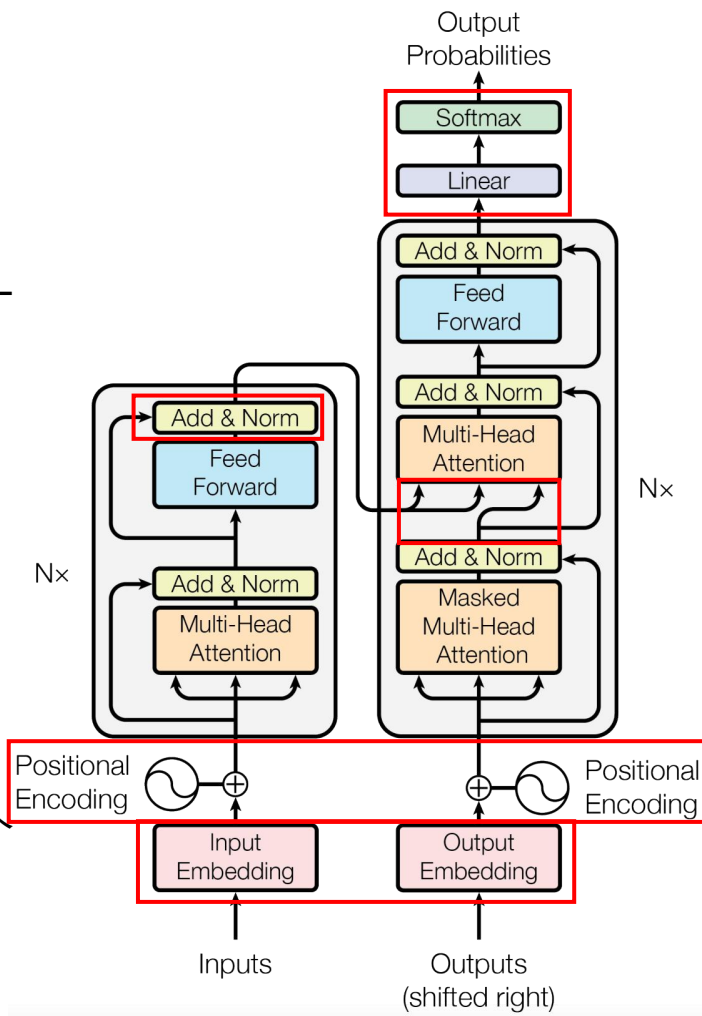


# Transformer Pytorch代码解读

# Transformer 架构

## 编码器

- 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<sub>1</sub> ~ Block<sub>N-1</sub>的输出: 输入到下个Block;
- Block<sub>N</sub>的输出: 输入到解码器的各层中。



## 解码器

- 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<sub>1</sub> ~ Block<sub>N-1</sub>的输出: 输入到下个Block;
- Block<sub>N</sub>的输出: 输入到后续的Linear层中。

# Transformer 架构

## Transformer代码架构

### 1. Transformer

#### Encoder

Word Embedding

PositionalEncoding

Encoder Layers (n\_layers = 6)

MultiHeadAttention (get\_attn\_pad\_mask)

ScaledDotProductAttention

Add (Residual Connection)

LayerNorm

PoswiseFeedForwardNet

Linear + ReLU + Linear

Add (Residual Connection)

LayerNorm

#### Decoder

Word Embedding

PositionalEncoding

Decoder Layers (n\_layers = 6)

MultiHeadAttention (get\_attn\_pad\_mask + get\_attn\_subsequence\_mask)

ScaledDotProductAttention

Add (Residual Connection)

LayerNorm

MultiHeadAttention (get\_attn\_pad\_mask)

ScaledDotProductAttention

Add (Residual Connection)

LayerNorm

PoswiseFeedForwardNet

Linear + ReLU + Linear

Add (Residual Connection)

LayerNorm

#### Projection

```
class Transformer(nn.Module):
    def __init__(self):
        super(Transformer, self).__init__()
        self.encoder = Encoder().cuda()
        self.decoder = Decoder().cuda()
        self.projection = nn.Linear(d_model, tgt_vocab_size, bias=False).cuda()

    def forward(self, enc_inputs, dec_inputs):
        """
        enc_inputs: [batch_size, src_len]
        dec_inputs: [batch_size, tgt_len]
        """
        # tensor to store decoder outputs
        # outputs = torch.zeros(batch_size, tgt_len, tgt_vocab_size).to(self.device)

        # enc_outputs: [batch_size, src_len, d_model], enc_self_attns: [n_layers, batch_size, n_heads, src_len, src_len]
        enc_outputs, enc_self_attns = self.encoder(enc_inputs)
        # dec_outpus: [batch_size, tgt_len, d_model], dec_self_attns: [n_layers, batch_size, n_heads, tgt_len, tgt_len],
        # dec_enc_attn: [n_layers, batch_size, tgt_len, src_len]
        dec_outputs, dec_self_attns, dec_enc_attns = self.decoder(dec_inputs, enc_inputs, enc_outputs)
        dec_logits = self.projection(dec_outputs) # dec_logits: [batch_size, tgt_len, tgt_vocab_size]
        return dec_logits.view(-1, dec_logits.size(-1)), enc_self_attns, dec_self_attns, dec_enc_attns
```

# Transformer 工作流程

## 编码过程

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

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

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

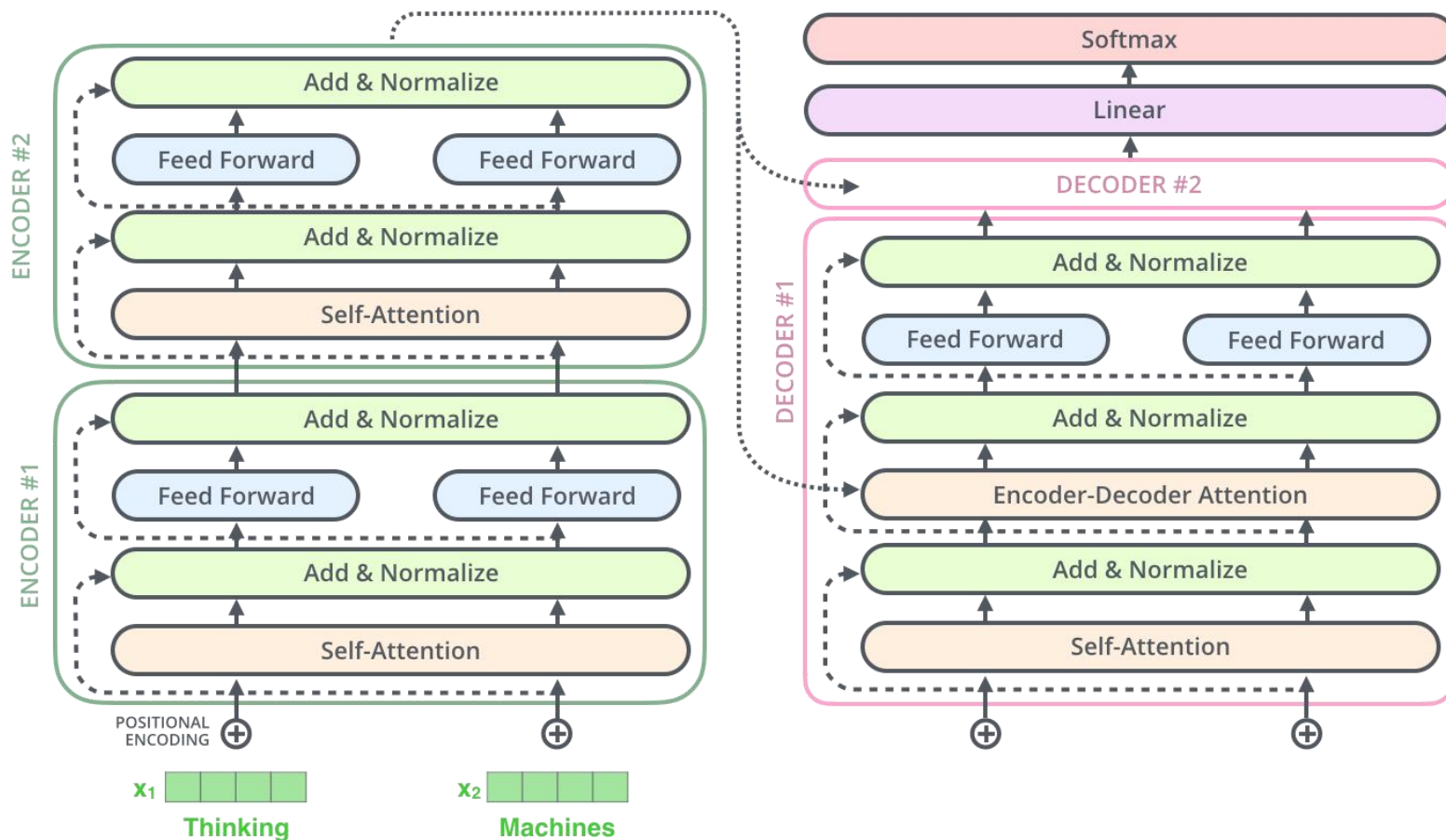
$$Q = \text{Linear}(X) = XW_Q$$

$$K = \text{Linear}(X) = XW_K$$

$$V = \text{Linear}(X) = XW_V$$

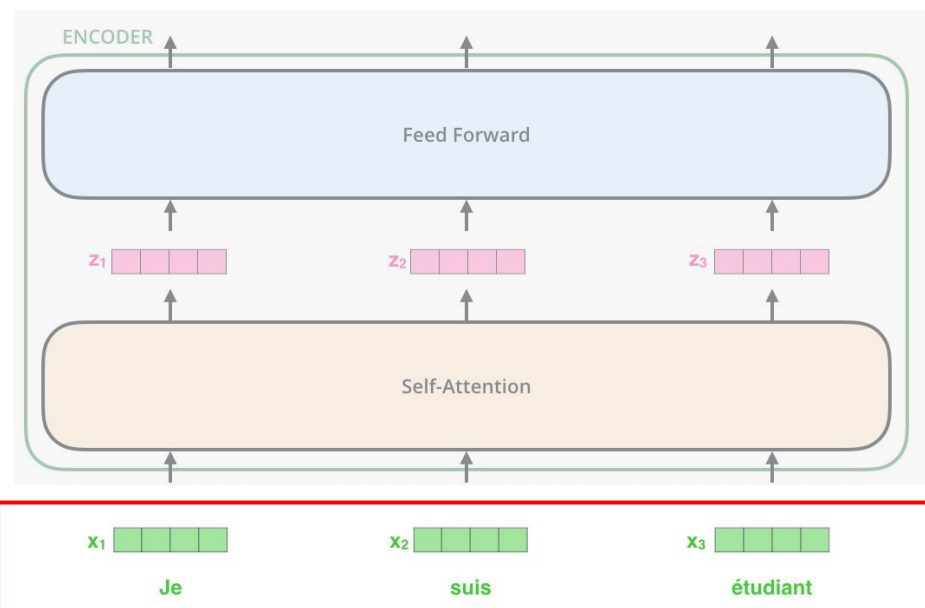
$$X_{\text{attention}} = \text{SelfAttention}(Q, K, V)$$

$$X = \text{Embedding Lookup}(X) + \text{Positional Encoding}$$



# Transformer 工作流程

## Word Embedding

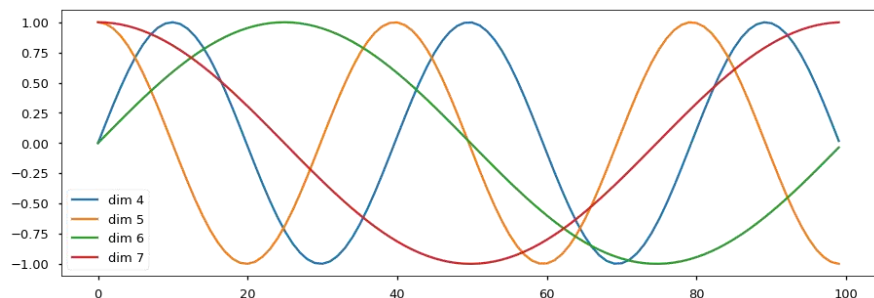
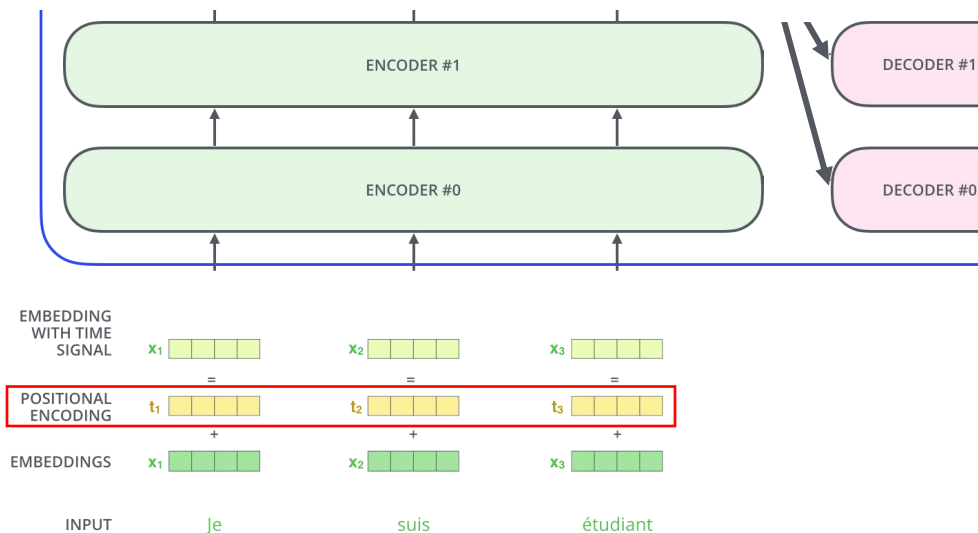


```
class Encoder(nn.Module):  
    def __init__(self):  
        super(Encoder, self).__init__()  
        self.src_emb = nn.Embedding(src_vocab_size, d_model)  
        self.pos_emb = PositionalEncoding(d_model)  
        self.layers = nn.ModuleList([EncoderLayer() for _ in range(n_layers)])
```

```
class Decoder(nn.Module):  
    def __init__(self):  
        super(Decoder, self).__init__()  
        self.tgt_emb = nn.Embedding(tgt_vocab_size, d_model)  
        self.pos_emb = PositionalEncoding(d_model)  
        self.layers = nn.ModuleList([DecoderLayer() for _ in range(n_layers)])
```

# Transformer 工作流程

## 位置编码 (Positional Encoding)



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

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

$$pos \in [0, \max\_sequence\_length)$$

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

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

```
class PositionalEncoding(nn.Module):
    def __init__(self, d_model, dropout=0.1, max_len=5000):
        super(PositionalEncoding, self).__init__()
        self.dropout = nn.Dropout(p=dropout)

        pe = torch.zeros(max_len, d_model)
        position = torch.arange(0, max_len, dtype=torch.float).unsqueeze(1)
        div_term = torch.exp(torch.arange(0, d_model, 2).float() * (-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).transpose(0, 1)
        self.register_buffer('pe', pe)

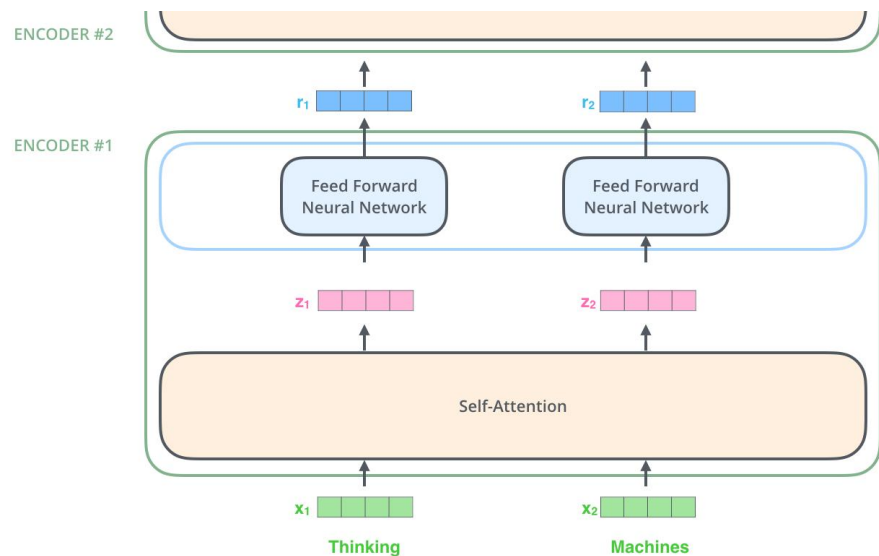
    def forward(self, x):
        ...
        x: [seq_len, batch_size, d_model]
        ...
        x = x + self.pe[:x.size(0), :]
        return self.dropout(x)
```

$$X = \text{Embedding Lookup}(X) + \text{Positional Encoding}$$



# Transformer 工作流程

## 编码过程 (Encoder)



```
class Encoder(nn.Module):
    def __init__(self):
        super(Encoder, self).__init__()
        self.src_emb = nn.Embedding(src_vocab_size, d_model)
        self.pos_emb = PositionalEncoding(d_model)
        self.layers = nn.ModuleList([EncoderLayer() for _ in range(n_layers)])

    def forward(self, enc_inputs):
        """
        enc_inputs: [batch_size, src_len]
        """
        enc_outputs = self.src_emb(enc_inputs) # [batch_size, src_len, d_model]
        enc_outputs = self.pos_emb(enc_outputs.transpose(0, 1)).transpose(0, 1) # [batch_size, src_len, d_model]
        enc_self_attn_mask = get_attn_pad_mask(enc_inputs, enc_inputs) # [batch_size, src_len, src_len]
        enc_self_attns = []
        for layer in self.layers:
            # enc_outputs: [batch_size, src_len, d_model], enc_self_attn: [batch_size, n_heads, src_len, src_len]
            enc_outputs, enc_self_attn = layer(enc_outputs, enc_self_attn_mask)
            enc_self_attns.append(enc_self_attn)
        return enc_outputs, enc_self_attns # 每个 block 有一个 attention mask
```

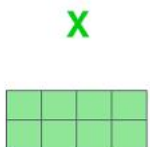
```
class EncoderLayer(nn.Module):
    def __init__(self):
        super(EncoderLayer, self).__init__()
        self.enc_self_attn = MultiHeadAttention()
        self.pos_ffn = PoswiseFeedForwardNet()

    def forward(self, enc_inputs, enc_self_attn_mask):
        """
        enc_inputs: [batch_size, src_len, d_model]
        enc_self_attn_mask: [batch_size, src_len, src_len]
        """
        # enc_outputs: [batch_size, src_len, d_model], attn: [batch_size, n_heads, src_len, src_len]
        enc_outputs, attn = self.enc_self_attn(enc_inputs, enc_inputs, enc_inputs, enc_self_attn_mask) # enc_inputs to same Q,K,V
        enc_outputs = self.pos_ffn(enc_outputs) # enc_outputs: [batch_size, src_len, d_model]
        return enc_outputs, attn
```

# Transformer 工作流程

## 编码过程 — Padding 操作

X: Thinking Machines

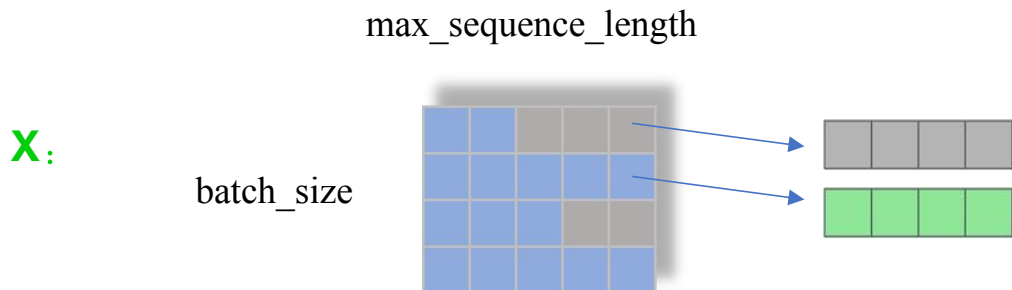


X 的维度: [sequence\_length, embedding\_dimension]



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]



```
class Encoder(nn.Module):
    def __init__(self):
        super(Encoder, self).__init__()
        self.src_emb = nn.Embedding(src_vocab_size, d_model)
        self.pos_emb = PositionalEncoding(d_model)
        self.layers = nn.ModuleList([EncoderLayer() for _ in range(n_layers)])

    def forward(self, enc_inputs):
        """
        enc_inputs: [batch_size, src_len]
        """
        enc_outputs = self.src_emb(enc_inputs) # [batch_size, src_len, d_model]
        enc_outputs = self.pos_emb(enc_outputs.transpose(0, 1)).transpose(0, 1) # [batch_size, src_len, d_model]
        enc_self_attn_mask = get_attn_pad_mask(enc_inputs, enc_inputs) # [batch_size, src_len, src_len]
        enc_self_attns = []
        for layer in self.layers:
            # enc_outputs: [batch_size, src_len, d_model], enc_self_attn: [batch_size, n_heads, src_len, src_len]
            enc_outputs, enc_self_attn = layer(enc_outputs, enc_self_attn_mask)
            enc_self_attns.append(enc_self_attn)
        return enc_outputs, enc_self_attns # 每个 block 有一个 attention mask
```

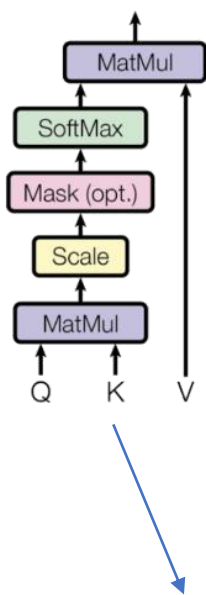
```
def get_attn_pad_mask(seq_q, seq_k):
    """
    seq_q: [batch_size, seq_len]
    seq_k: [batch_size, seq_len]
    seq_len could be src_len or it could be tgt_len
    seq_len in seq_q and seq_len in seq_k maybe not equal
    """
    batch_size, len_q = seq_q.size()
    batch_size, len_k = seq_k.size()
    # eq(zero) is PAD token
    pad_attn_mask = seq_k.data.eq(0).unsqueeze(1) # [batch_size, 1, len_k], False is masked
    return pad_attn_mask.expand(batch_size, len_q, len_k) # [batch_size, len_q, len_k]
```



# Transformer 工作流程

## 编码过程 — Encoder Multi-Head Self-Attention

Scaled dot-product attention



$$\begin{aligned} Q &= \text{Linear}(X) = XW_Q \\ K &= \text{Linear}(X) = XW_K \\ V &= \text{Linear}(X) = XW_V \end{aligned}$$

第一步：生成 Q、K、V，辅助计算注意力机制

```
class EncoderLayer(nn.Module):
    def __init__(self):
        super(EncoderLayer, self).__init__()
        self.enc_self_attn = MultiHeadAttention()
        self.pos_ffn = PoswiseFeedForwardNet()

    def forward(self, enc_inputs, enc_self_attn_mask):
        """
        enc_inputs: [batch_size, src_len, d_model]
        enc_self_attn_mask: [batch_size, src_len, src_len]
        """
        # enc_outputs: [batch_size, src_len, d_model], attn: [batch_size, n_heads, src_len, src_len]
        enc_outputs, attn = self.enc_self_attn(enc_inputs, enc_inputs, enc_inputs, enc_self_attn_mask) # enc_inputs to same Q,K,V
        enc_outputs = self.pos_ffn(enc_outputs) # enc_outputs: [batch_size, src_len, d_model]
        return enc_outputs, attn
```

```
class MultiHeadAttention(nn.Module):
    def __init__(self):
        super(MultiHeadAttention, self).__init__()
        self.W_Q = nn.Linear(d_model, d_k * n_heads, bias=False)
        self.W_K = nn.Linear(d_model, d_k * n_heads, bias=False)
        self.W_V = nn.Linear(d_model, d_v * n_heads, bias=False)
        self.fc = nn.Linear(n_heads * d_v, d_model, bias=False)

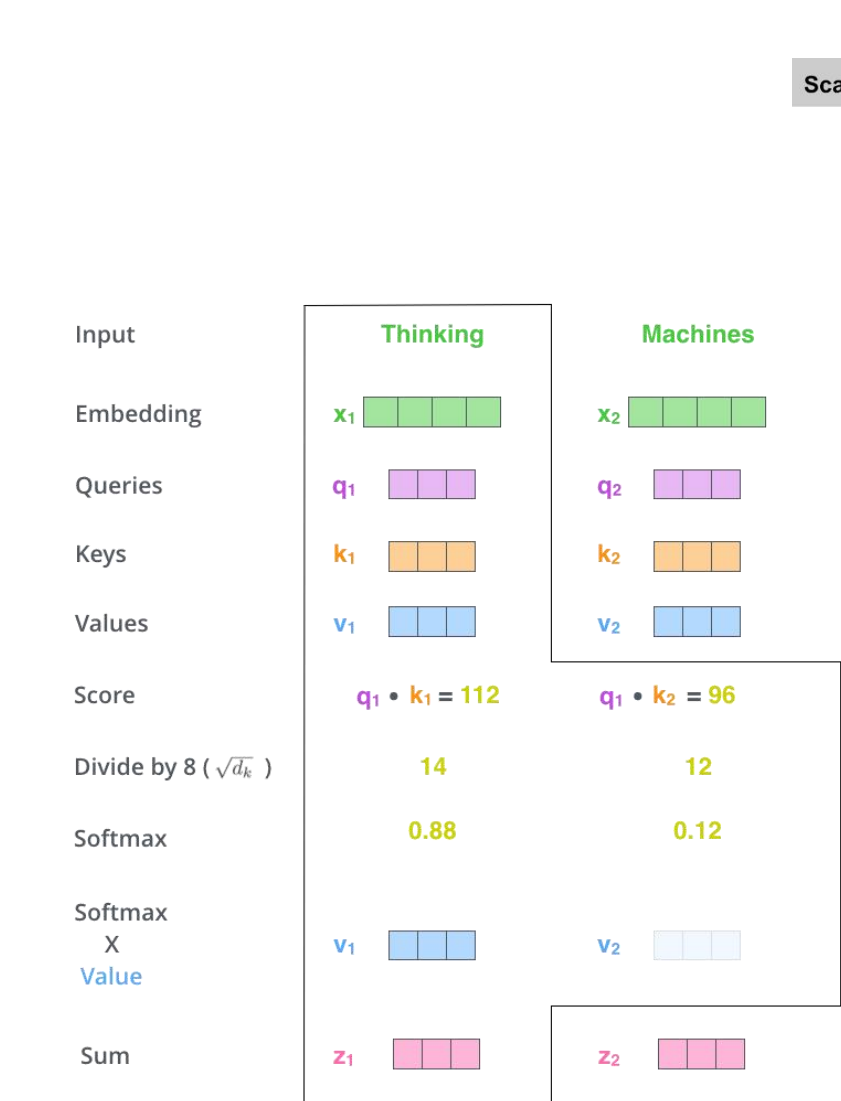
    def forward(self, input_Q, input_K, input_V, attn_mask):
        """
        input_Q: [batch_size, len_q, d_model]
        input_K: [batch_size, len_k, d_model]
        input_V: [batch_size, len_v(=len_k), d_model]
        attn_mask: [batch_size, seq_len, seq_len]
        """
        residual, batch_size = input_Q, input_Q.size(0)
        # (B, S, D) -proj-> (B, S, D_new) -split-> (B, S, H, W) -trans-> (B, H, S, W)
        Q = self.W_Q(input_Q).view(batch_size, -1, n_heads, d_k).transpose(1,2) # Q: [batch_size, n_heads, len_q, d_k]
        K = self.W_K(input_K).view(batch_size, -1, n_heads, d_k).transpose(1,2) # K: [batch_size, n_heads, len_k, d_k]
        V = self.W_V(input_V).view(batch_size, -1, n_heads, d_v).transpose(1,2) # V: [batch_size, n_heads, len_v(=len_k), d_v]

        attn_mask = attn_mask.unsqueeze(1).repeat(1, n_heads, 1, 1) # attn_mask : [batch_size, n_heads, seq_len, seq_len]

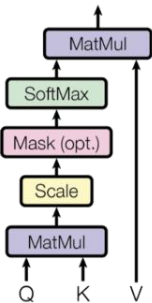
        # context: [batch_size, n_heads, len_q, d_v], attn: [batch_size, n_heads, len_q, len_k]
        context, attn = ScaledDotProductAttention()(Q, K, V, attn_mask)
        context = context.transpose(1, 2).reshape(batch_size, -1, n_heads * d_v) # context: [batch_size, len_q, n_heads * d_v]
        output = self.fc(context) # [batch_size, len_q, d_model]
        return nn.LayerNorm(d_model).cuda()(output + residual), attn
```

# Transformer 工作流程

## 编码过程 — Scaled Dot Product Attention



Scaled dot-product attention



核心公式:  $Attention(Q, K, V) = softmax(\frac{QK^T}{\sqrt{d_k}}) V$

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

$e^0 = 1$

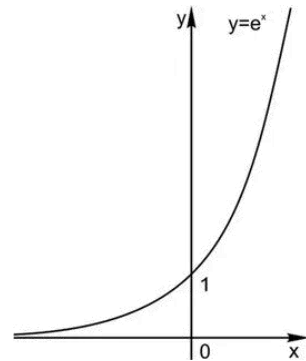
$e^{-\infty} \rightarrow 0$

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



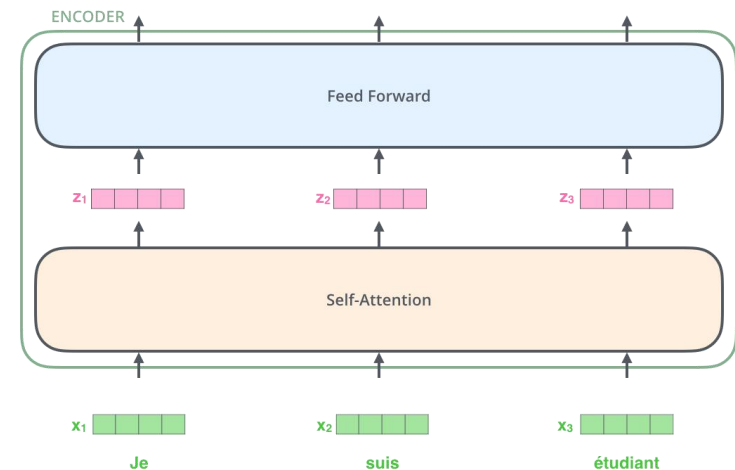
```
class ScaledDotProductAttention(nn.Module):
    def __init__(self):
        super(ScaledDotProductAttention, self).__init__()

    def forward(self, Q, K, V, attn_mask):
        """
        Q: [batch_size, n_heads, len_q, d_k]
        K: [batch_size, n_heads, len_k, d_k]
        V: [batch_size, n_heads, len_v(=len_k), d_v]
        attn_mask: [batch_size, n_heads, seq_len, seq_len]
        """
        scores = torch.matmul(Q, K.transpose(-1, -2)) / np.sqrt(d_k) # scores : [batch_size, n_heads, len_q, len_k]
        scores.masked_fill_(attn_mask, -1e9) # Fills elements of self tensor with value where mask is True.

        attn = nn.Softmax(dim=-1)(scores)
        context = torch.matmul(attn, V) # [batch_size, n_heads, len_q, d_v]
        return context, attn
```

# Transformer 工作流程

## 编码过程 — Encoder Multi-Head Self-Attention



1) Concatenate all the attention heads



2) Multiply with a weight matrix  $W^o$  that was trained jointly with the model

$\times$

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



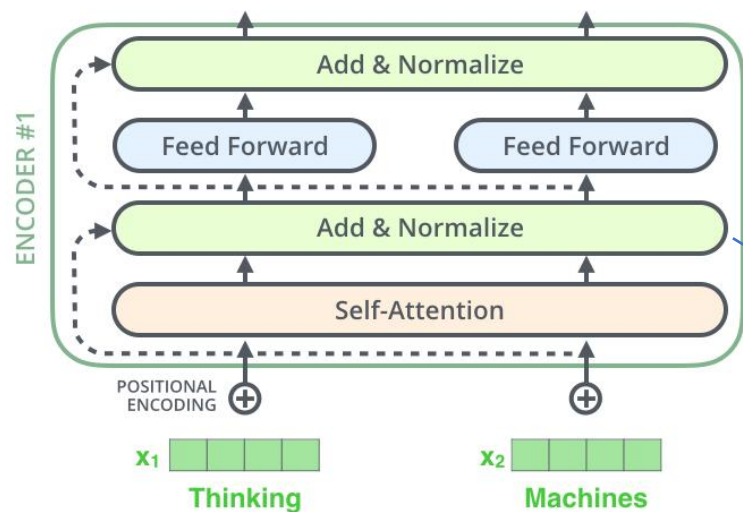
```
class MultiHeadAttention(nn.Module):
    def __init__(self):
        super(MultiHeadAttention, self).__init__()
        self.W_Q = nn.Linear(d_model, d_k * n_heads, bias=False)
        self.W_K = nn.Linear(d_model, d_k * n_heads, bias=False)
        self.W_V = nn.Linear(d_model, d_v * n_heads, bias=False)
        self.fc = nn.Linear(n_heads * d_v, d_model, bias=False)
    def forward(self, input_Q, input_K, input_V, attn_mask):
        ...
        input_Q: [batch_size, len_q, d_model]
        input_K: [batch_size, len_k, d_model]
        input_V: [batch_size, len_v(=len_k), d_model]
        attn_mask: [batch_size, seq_len, seq_len]
        ...
        residual, batch_size = input_Q, input_Q.size(0)
        # (B, S, D) -proj-> (B, S, D_new) -split-> (B, S, H, W) -trans-> (B, H, S, W)
        Q = self.W_Q(input_Q).view(batch_size, -1, n_heads, d_k).transpose(1,2) # Q: [batch_size, n_heads, len_q, d_k]
        K = self.W_K(input_K).view(batch_size, -1, n_heads, d_k).transpose(1,2) # K: [batch_size, n_heads, len_k, d_k]
        V = self.W_V(input_V).view(batch_size, -1, n_heads, d_v).transpose(1,2) # V: [batch_size, n_heads, len_v(=len_k), d_v]

        attn_mask = attn_mask.unsqueeze(1).repeat(1, n_heads, 1, 1) # attn_mask : [batch_size, n_heads, seq_len, seq_len]

        # context: [batch_size, n_heads, len_q, d_v], attn: [batch_size, n_heads, len_q, len_k]
        context, attn = ScaledDotProductAttention()(Q, K, V, attn_mask)
        context = context.transpose(1, 2).reshape(batch_size, -1, n_heads * d_v) # context: [batch_size, len_q, n_heads * d_v]
        output = self.fc(context) # [batch_size, len_q, d_model]
        return nn.LayerNorm(d_model).cuda()(output + residual), attn
```

# Transformer 工作流程

## 编码过程 — Add & Normalize



$$\begin{aligned} X_{\text{hidden}} &= X_{\text{attention}} + X_{\text{hidden}} \\ X_{\text{hidden}} &= \text{LayerNorm}(X_{\text{hidden}}) \end{aligned}$$

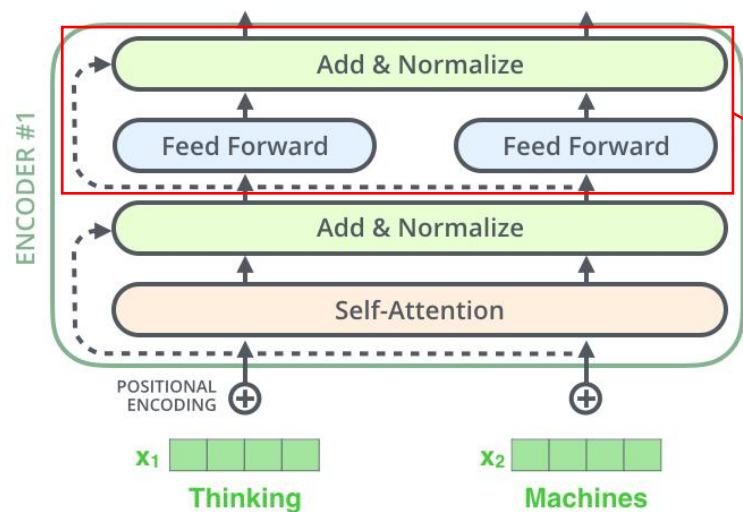
```
class MultiHeadAttention(nn.Module):
    def __init__(self):
        super(MultiHeadAttention, self).__init__()
        self.W_Q = nn.Linear(d_model, d_k * n_heads, bias=False)
        self.W_K = nn.Linear(d_model, d_k * n_heads, bias=False)
        self.W_V = nn.Linear(d_model, d_v * n_heads, bias=False)
        self.fc = nn.Linear(n_heads * d_v, d_model, bias=False)
    def forward(self, input_Q, input_K, input_V, attn_mask):
        ...
        input_Q: [batch_size, len_q, d_model]
        input_K: [batch_size, len_k, d_model]
        input_V: [batch_size, len_v(=len_k), d_model]
        attn_mask: [batch_size, seq_len, seq_len]
        ...
        residual, batch_size = input_Q, input_Q.size(0)
        # (B, S, D) ->proj-> (B, S, D_new) ->split-> (B, S, H, W) ->trans-> (B, H, S, W)
        Q = self.W_Q(input_Q).view(batch_size, -1, n_heads, d_k).transpose(1,2) # Q: [batch_size, n_heads, len_q, d_k]
        K = self.W_K(input_K).view(batch_size, -1, n_heads, d_k).transpose(1,2) # K: [batch_size, n_heads, len_k, d_k]
        V = self.W_V(input_V).view(batch_size, -1, n_heads, d_v).transpose(1,2) # V: [batch_size, n_heads, len_v(=len_k), d_v]

        attn_mask = attn_mask.unsqueeze(1).repeat(1, n_heads, 1, 1) # attn_mask : [batch_size, n_heads, seq_len, seq_len]

        # context: [batch_size, n_heads, len_q, d_v], attn: [batch_size, n_heads, len_q, len_k]
        context, attn = ScaledDotProductAttention()(Q, K, V, attn_mask)
        context = context.transpose(1, 2).reshape(batch_size, -1, n_heads * d_v) # context: [batch_size, len_q, n_heads * d_v]
        output = self.fc(context) # [batch_size, len_q, d_model]
        return nn.LayerNorm(d_model).cuda()(output + residual), attn
```

# Transformer 工作流程

## 编码过程 — Feed Forward Net



```
class EncoderLayer(nn.Module):
    def __init__(self):
        super(EncoderLayer, self).__init__()
        self.enc_self_attn = MultiHeadAttention()
        self.pos_ffn = PoswiseFeedForwardNet()

    def forward(self, enc_inputs, enc_self_attn_mask):
        ...
        enc_inputs: [batch_size, src_len, d_model]
        enc_self_attn_mask: [batch_size, src_len, src_len]
        ...
        # enc_outputs: [batch_size, src_len, d_model], attn: [batch_size, n_heads, src_len, src_len]
        enc_outputs, attn = self.enc_self_attn(enc_inputs, enc_inputs, enc_inputs, enc_self_attn_mask) # enc_inputs to same Q,K,V
        enc_outputs = self.pos_ffn(enc_outputs) # enc_outputs: [batch_size, src_len, d_model]
        return enc_outputs, attn
```

```
class PoswiseFeedForwardNet(nn.Module):
    def __init__(self):
        super(PoswiseFeedForwardNet, self).__init__()
        self.fc = nn.Sequential(
            nn.Linear(d_model, d_ff, bias=False),
            nn.ReLU(),
            nn.Linear(d_ff, d_model, bias=False)
        )
        self.ln = nn.LayerNorm(d_model).cuda()
    def forward(self, inputs):
        ...
        这里进行了残差连接: Add&LayerNorm
        inputs: [batch_size, seq_len, d_model]
        ...
        residual = inputs
        output = self.fc(inputs)
        return self.ln(output + residual) # [batch_size, seq_len, d_model]
```

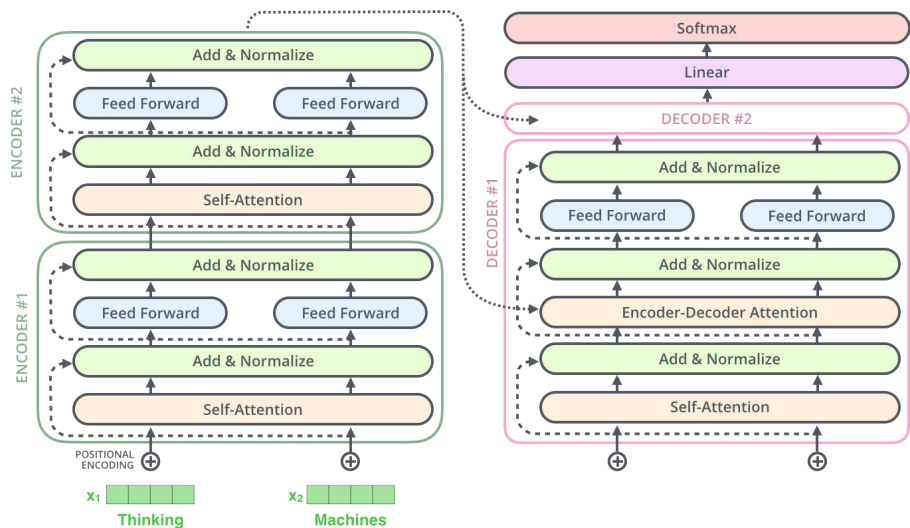
$$X_{\text{hidden}} = \text{Linear}(\text{ReLU}(\text{Linear}(X_{\text{attention}})))$$

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



# Transformer 工作流程

## 解码过程



```
class Decoder(nn.Module):
    def __init__(self):
        super(Decoder, self).__init__()
        self.tgt_emb = nn.Embedding(tgt_vocab_size, d_model)
        self.pos_emb = PositionalEncoding(d_model)
        self.layers = nn.ModuleList([DecoderLayer() for _ in range(n_layers)])

    def forward(self, dec_inputs, enc_inputs, enc_outputs):
        """
        dec_inputs: [batch_size, tgt_len]
        enc_inputs: [batch_size, src_len]
        enc_outputs: [batch_size, src_len, d_model]
        """
        dec_outputs = self.tgt_emb(dec_inputs) # [batch_size, tgt_len, d_model]
        dec_outputs = self.pos_emb(dec_outputs.transpose(0, 1)).transpose(0, 1).cuda() # [batch_size, tgt_len, d_model]
        dec_self_attn_pad_mask = get_attn_pad_mask(dec_inputs, dec_inputs).cuda() # [batch_size, tgt_len, tgt_len]
        dec_self_attn_subsequence_mask = get_attn_subsequence_mask(dec_inputs).cuda() # [batch_size, tgt_len, tgt_len]
        dec_self_attn_mask = torch.gt((dec_self_attn_pad_mask + dec_self_attn_subsequence_mask), 0).cuda() # [batch_size, tgt_len, tgt_len] gt 函数: greater than

        dec_enc_attn_mask = get_attn_pad_mask(dec_inputs, enc_inputs) # [batch_size, tgt_len, src_len]

        dec_self_attns, dec_enc_attns = [], []
        for layer in self.layers:
            # dec_outputs: [batch_size, tgt_len, d_model], dec_self_attn: [batch_size, n_heads, tgt_len, tgt_len], dec_enc_attn: [batch_size, h_heads, tgt_len, src_len]
            dec_outputs, dec_self_attn, dec_enc_attn = layer(dec_outputs, enc_outputs, dec_self_attn_mask, dec_enc_attn_mask)
            dec_self_attns.append(dec_self_attn)
            dec_enc_attns.append(dec_enc_attn)
        return dec_outputs, dec_self_attns, dec_enc_attns
```

```
class DecoderLayer(nn.Module):
    def __init__(self):
        super(DecoderLayer, self).__init__()
        self.dec_self_attn = MultiHeadAttention()
        self.dec_enc_attn = MultiHeadAttention()
        self.pos_ffn = PoswiseFeedForwardNet()

    def forward(self, dec_inputs, enc_outputs, dec_self_attn_mask, dec_enc_attn_mask):
        """
        dec_inputs: [batch_size, tgt_len, d_model]
        enc_outputs: [batch_size, src_len, d_model]
        dec_self_attn_mask: [batch_size, tgt_len, tgt_len]
        dec_enc_attn_mask: [batch_size, tgt_len, src_len]
        """
        # dec_outputs: [batch_size, tgt_len, d_model], dec_self_attn: [batch_size, n_heads, tgt_len, tgt_len]
        dec_outputs, dec_self_attn = self.dec_self_attn(dec_inputs, dec_inputs, dec_inputs, dec_self_attn_mask)

        # dec_outputs: [batch_size, tgt_len, d_model], dec_enc_attn: [batch_size, h_heads, tgt_len, src_len]
        dec_outputs, dec_enc_attn = self.dec_enc_attn(dec_outputs, enc_outputs, enc_outputs, dec_enc_attn_mask)

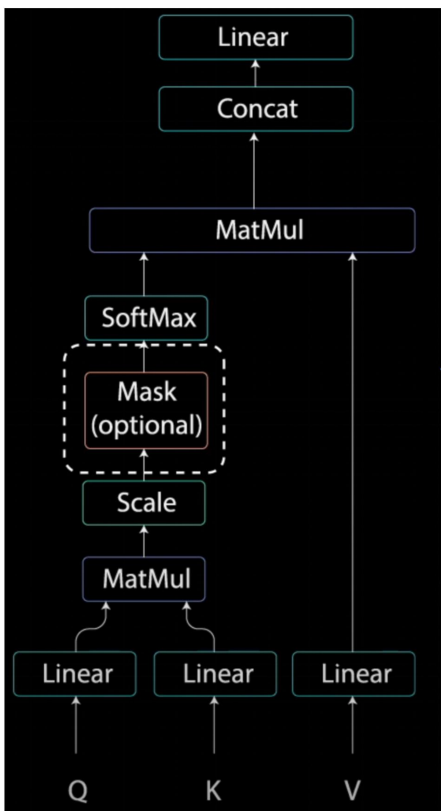
        dec_outputs = self.pos_ffn(dec_outputs) # [batch_size, tgt_len, d_model]

        return dec_outputs, dec_self_attn, dec_enc_attn
```



# Transformer 工作流程

## 解码过程 — Decoder Masked Self-Attention



```
def get_attn_subsequence_mask(seq):  
    ...  
    seq: [batch_size, tgt_len]  
    ...  
    attn_shape = [seq.size(0), seq.size(1), seq.size(1)]  
    subsequence_mask = np.triu(np.ones(attn_shape), k=1) # Upper triangular matrix  
    subsequence_mask = torch.from_numpy(subsequence_mask).byte()  
    return subsequence_mask # [batch_size, tgt_len, tgt_len]
```

Scaled Scores					Look-Ahead Mask					Masked Scores				
	<start>	I	am	fine										
<start>	0.7	0.1	0.1	0.1		0	-inf	-inf	-inf		0.7	-inf	-inf	-inf
I	0.1	0.6	0.2	0.1		0	0	-inf	-inf		0.1	0.6	-inf	-inf
am	0.1	0.2	0.6	0.1		0	0	0	-inf		0.1	0.2	0.6	-inf
fine	0.1	0.3	0.3	0.3		0	0	0	0		0.1	0.3	0.3	0.3

	<start>	I	am	fine
<start>	1	0	0	0
I	0.37	0.62	0	0
am	0.26	0.31	0.43	0
fine	0.21	0.26	0.26	0.26

Softmax(

	<start>	I	am	fine
<start>	0.7	-inf	-inf	-inf
I	0.1	0.6	-inf	-inf
am	0.1	0.2	0.6	-inf
fine	0.1	0.3	0.3	0.3

) =

# Transformer 工作流程

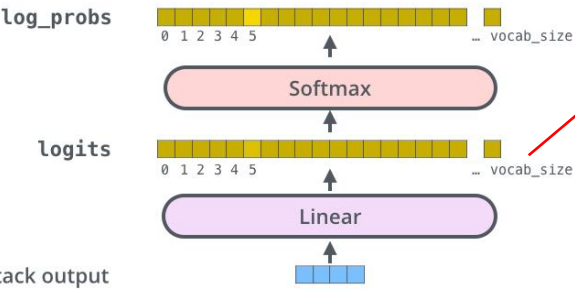
## 解码过程 —— Linear & Softmax

Which word in our vocabulary is associated with this index?

am

Get the index of the cell with the highest value (argmax)

5



```
class Transformer(nn.Module):
    def __init__(self):
        super(Transformer, self).__init__()
        self.encoder = Encoder().cuda()
        self.decoder = Decoder().cuda()
        self.projection = nn.Linear(d_model, tgt_vocab_size, bias=False).cuda()
    def forward(self, enc_inputs, dec_inputs):
        """
        enc_inputs: [batch_size, src_len]
        dec_inputs: [batch_size, tgt_len]
        """
        # tensor to store decoder outputs
        # outputs = torch.zeros(batch_size, tgt_len, tgt_vocab_size).to(self.device)

        # enc_outputs: [batch_size, src_len, d_model], enc_self_attns: [n_layers, batch_size, n_heads, src_len, src_len]
        enc_outputs, enc_self_attns = self.encoder(enc_inputs)
        # dec_outpus: [batch_size, tgt_len, d_model], dec_self_attns: [n_layers, batch_size, n_heads, tgt_len, tgt_len], dec_enc_attn: [n_layers, batch_size, tgt_len, src_len]
        dec_outputs, dec_self_attns, dec_enc_attns = self.decoder(dec_inputs, enc_inputs, enc_outputs)
        dec_logits = self.projection(dec_outputs) # dec_logits: [batch_size, tgt_len, tgt_vocab_size]
        return dec_logits.view(-1, dec_logits.size(-1)), enc_self_attns, dec_self_attns, dec_enc_attns
```