

```

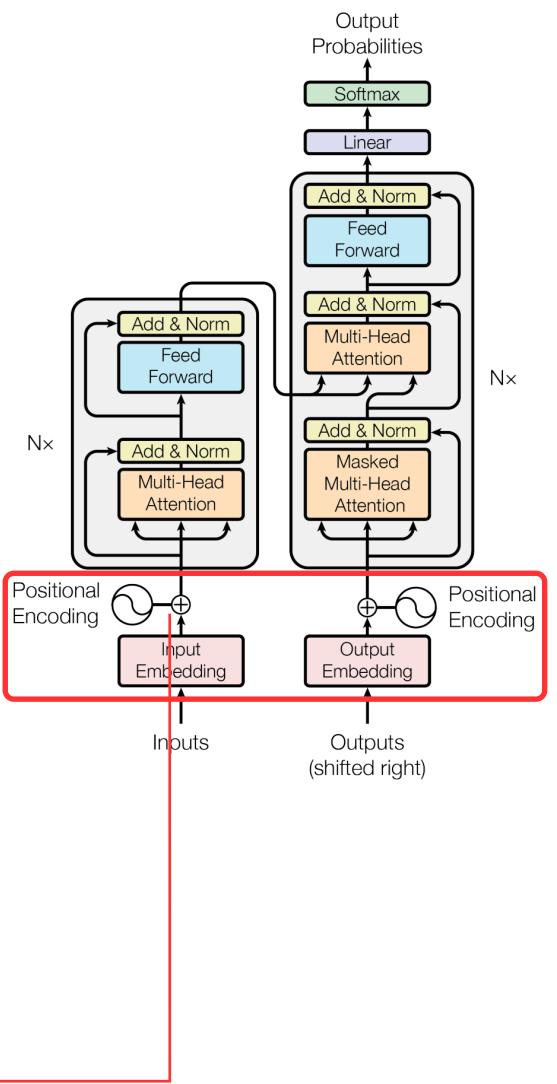
1 class Embeddings(nn.Module):
2     def __init__(self, d_model, vocab):
3         super(Embeddings, self).__init__()
4         self.lut = nn.Embedding(vocab, d_model)
5         self.d_model = d_model
6
7     def forward(self, x):
8         return self.lut(x) * math.sqrt(self.d_model)

```

```

1 class PositionalEncoding(nn.Module):
2     "Implement the PE function."
3
4     def __init__(self, d_model, dropout, max_len=5000):
5         super(PositionalEncoding, self).__init__()
6         self.dropout = nn.Dropout(p=dropout)
7
8         # Compute the positional encodings once in log space.
9         pe = torch.zeros(max_len, d_model)
10        position = torch.arange(0, max_len).unsqueeze(1)
11        div_term = torch.exp(
12            torch.arange(0, d_model, 2) * -(math.log(10000.0) / d_model)
13        )
14        pe[:, 0::2] = torch.sin(position * div_term)
15        pe[:, 1::2] = torch.cos(position * div_term)
16        pe = pe.unsqueeze(0)
17        self.register_buffer("pe", pe)
18
19    def forward(self, x):
20        x = x + self.pe[:, :x.size(1)].requires_grad_(False) ←
21        return self.dropout(x)

```



```

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

```



```

1 def attention(query, key, value, mask=None, dropout=None):
2     "Compute 'Scaled Dot Product Attention'"      这个是交换最后两个维
3     d_k = query.size(-1)                         度, 理由是对齐
4     scores = torch.matmul(query, key.transpose(-2, -1)) / math.sqrt(d_k)
5     if mask is not None:
6         scores = scores.masked_fill(mask == 0, -1e9)
7     p_attn = scores.softmax(dim=-1)
8     if dropout is not None:
9         p_attn = dropout(p_attn)
10    return torch.matmul(p_attn, value), p_attn

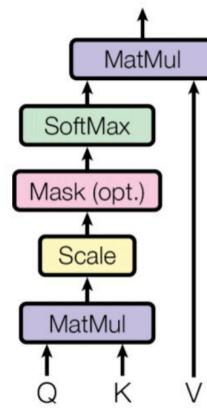
```



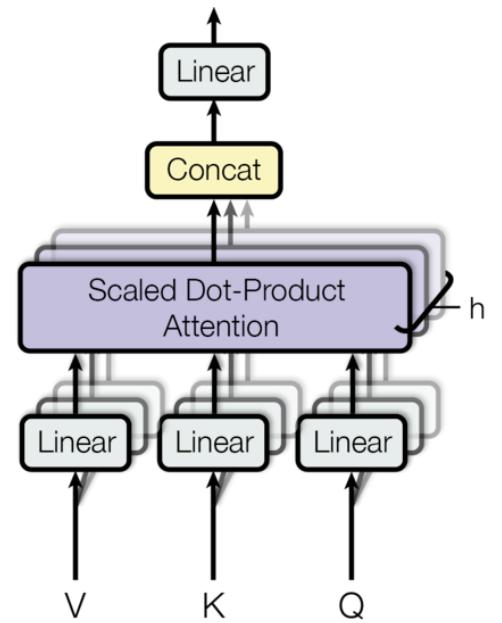
```

1 class MultiHeadedAttention(nn.Module):
2     def __init__(self, h, d_model, dropout=0.1):
3         "Take in model size and number of heads."
4         super(MultiHeadedAttention, self).__init__()
5         assert d_model % h == 0
6         # We assume d_v always equals d_k
7         self.d_k = d_model // h
8         self.h = h                                复制4个不一样的linear layer
9         self.linears = clones(nn.Linear(d_model, d_model), 4)
10        self.attn = None
11        self.dropout = nn.Dropout(p=dropout)
12
13    def forward(self, query, key, value, mask=None):
14        "Implements Figure 2"
15        if mask is not None:
16            # Same mask applied to all h heads.
17            mask = mask.unsqueeze(1)
18            nbatches = query.size(0)
19
20            # 1) Do all the linear projections in batch from d_model => h x d_k
21            query, key, value = [          这边用掉3个linear layer
22                lin(x).view(nbatches, -1, self.h, self.d_k).transpose(1, 2)
23                for lin, x in zip(self.linears, (query, key, value))
24            ]
25
26            # 2) Apply attention on all the projected vectors in batch.
27            x, self.attn = attention(
28                query, key, value, mask=mask, dropout=self.dropout
29            )
30
31            # 3) "Concat" using a view and apply a final linear.
32            x = (
33                x.transpose(1, 2)
34                .contiguous()           这边就是 concat 部分
35                .view(nbatches, -1, self.h * self.d_k)
36            )
37            del query
38            del key
39            del value   出来的时候用掉最后一个 linear layer
40            return self.linears[-1](x)

```



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





### Feed Forward

```

1 class PositionwiseFeedForward(nn.Module):
2     "Implements FFN equation."
3
4     def __init__(self, d_model, d_ff, dropout=0.1):
5         super(PositionwiseFeedForward, self).__init__()
6         self.w_1 = nn.Linear(d_model, d_ff)
7         self.w_2 = nn.Linear(d_ff, d_model)
8         self.dropout = nn.Dropout(dropout)
9
10    def forward(self, x):
11        return self.w_2(self.dropout(self.w_1(x).relu()))

```

$$FFN(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

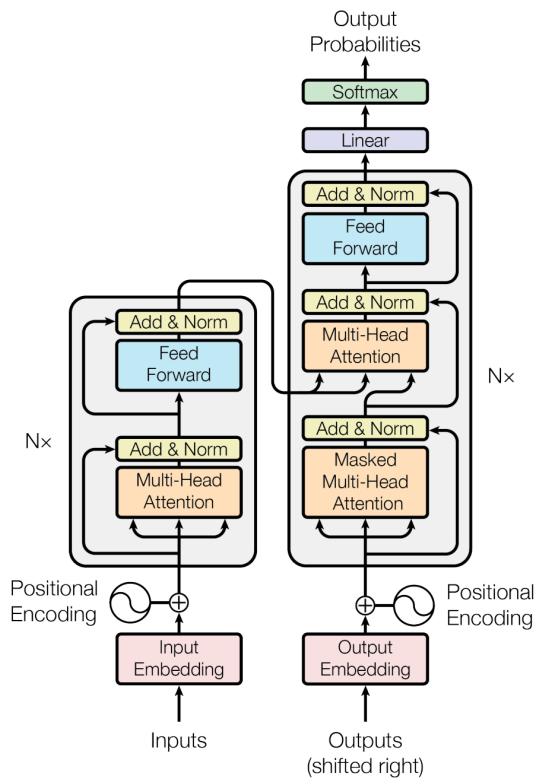


### Norm

```

1 class LayerNorm(nn.Module):
2     "Construct a layernorm module (See citation for details)."
3
4     def __init__(self, features, eps=1e-6):
5         super(LayerNorm, self).__init__()
6         self.a_2 = nn.Parameter(torch.ones(features))
7         self.b_2 = nn.Parameter(torch.zeros(features))
8         self.eps = eps
9
10    def forward(self, x):
11        mean = x.mean(-1, keepdim=True)
12        std = x.std(-1, keepdim=True)
13        return self.a_2 * (x - mean) / (std + self.eps) + self.b_2

```



### 最后两层

```

1 class Generator(nn.Module):
2     "Define standard linear + softmax generation step."
3
4     def __init__(self, d_model, vocab):
5         super(Generator, self).__init__()
6         self.proj = nn.Linear(d_model, vocab)
7
8     def forward(self, x):
9         return log_softmax(self.proj(x), dim=-1)

```

```

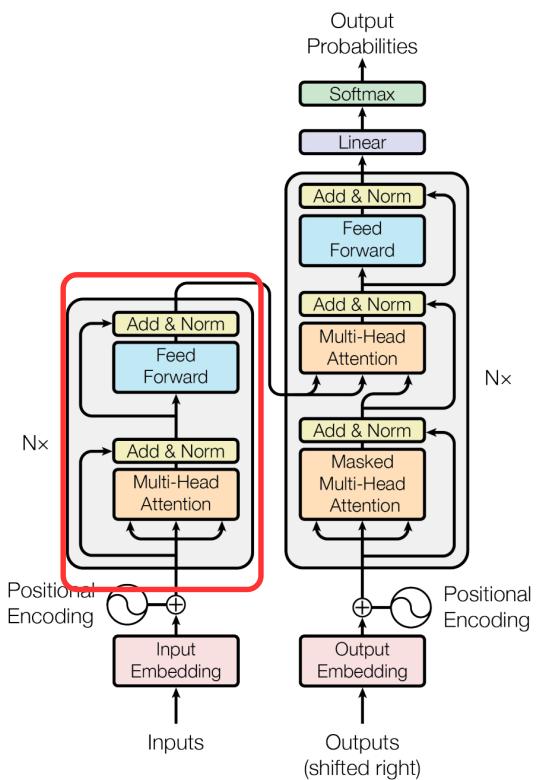
1 def clones(module, N):
2     "Produce N identical layers."
3     return nn.ModuleList([copy.deepcopy(module) for _ in range(N)])

```

```

1 class SublayerConnection(nn.Module):
2     """
3         A residual connection followed by a layer norm.
4         Note for code simplicity the norm is first as opposed to last.
5     """
6
7     def __init__(self, size, dropout):
8         super(SublayerConnection, self).__init__()
9         self.norm = LayerNorm(size)
10        self.dropout = nn.Dropout(dropout)
11
12    def forward(self, x, sublayer):
13        "Apply residual connection to any sublayer with the same size."
14        return x + self.dropout(sublayer(self.norm(x)))

```



```

1 class EncoderLayer(nn.Module):
2     "Encoder is made up of self-attn and feed forward (defined below)"
3
4     def __init__(self, size, self_attn, feed_forward, dropout):
5         super(EncoderLayer, self).__init__()
6         self.self_attn = self_attn
7         self.feed_forward = feed_forward
8         self.sublayer = clones(SublayerConnection(size, dropout), 2)
9         self.size = size
10
11    def forward(self, x, mask):
12        "Follow Figure 1 (left) for connections."
13        x = self.sublayer[0](x, lambda x: self.self_attn(x, x, x, mask))
14        return self.sublayer[1](x, self.feed_forward)

```

这里是一个 encoder

这个就是有几个encoder了  
注意：每一个encoder之后都要 norm

```

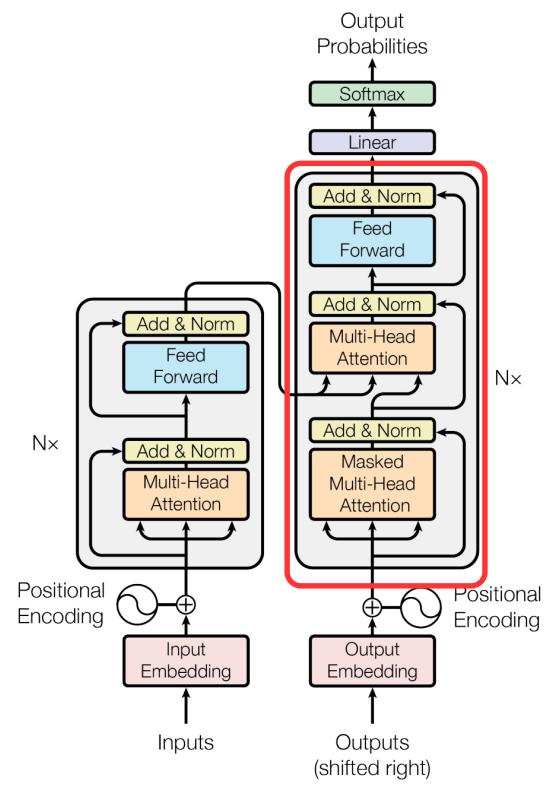
1 class Encoder(nn.Module):
2     "Core encoder is a stack of N layers"
3
4     def __init__(self, layer, N):
5         super(Encoder, self).__init__()
6         self.layers = clones(layer, N)
7         self.norm = LayerNorm(layer.size)
8
9     def forward(self, x, mask):
10        "Pass the input (and mask) through each layer in turn."
11        for layer in self.layers:
12            x = layer(x, mask)
13        return self.norm(x)

```

```

1 class DecoderLayer(nn.Module):
2     "Decoder is made of self-attn, src-attn, and feed forward (defined
3      below)"
4
5     def __init__(self, size, self_attn, src_attn, feed_forward, dropout):
6         super(DecoderLayer, self).__init__()
7         self.size = size
8         self.self_attn = self_attn
9         self.src_attn = src_attn
10        self.feed_forward = feed_forward
11        self.sublayer = clones(SublayerConnection(size, dropout), 3)
12
13    def forward(self, x, memory, src_mask, tgt_mask):
14        "Follow Figure 1 (right) for connections."
15        m = memory                                         memory是从encoder出来的
16        x = self.sublayer[0](x, lambda x: self.self_attn(x, x, x, tgt_mask))
17        x = self.sublayer[1](x, lambda x: self.src_attn(x, m, m, src_mask))
18        return self.sublayer[2](x, self.feed_forward)

```



```

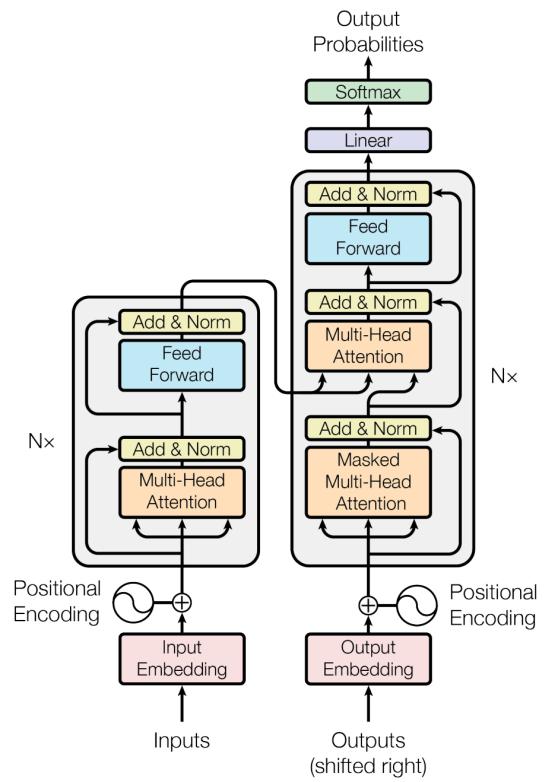
1 class Decoder(nn.Module):
2     "Generic N layer decoder with masking."
3
4     def __init__(self, layer, N):
5         super(Decoder, self).__init__()
6         self.layers = clones(layer, N)
7         self.norm = LayerNorm(layer.size)
8
9     def forward(self, x, memory, src_mask, tgt_mask):
10        for layer in self.layers:
11            x = layer(x, memory, src_mask, tgt_mask)
12        return self.norm(x)

```

```

1 class EncoderDecoder(nn.Module):
2     """
3         A standard Encoder-Decoder architecture. Base for this and many
4         other models.
5     """
6
7     def __init__(self, encoder, decoder, src_embed, tgt_embed, generator):
8         super(EncoderDecoder, self).__init__()
9         self.encoder = encoder
10        self.decoder = decoder
11        self.src_embed = src_embed
12        self.tgt_embed = tgt_embed
13        self.generator = generator
14
15    def forward(self, src, tgt, src_mask, tgt_mask):
16        "Take in and process masked src and target sequences."
17        return self.decode(self.encode(src, src_mask), src_mask, tgt, tgt_mask)
18
19    def encode(self, src, src_mask):
20        return self.encoder(self.src_embed(src), src_mask)
21
22    def decode(self, memory, src_mask, tgt, tgt_mask):
23        return self.decoder(self.tgt_embed(tgt), memory, src_mask, tgt_mask)

```



```

1 def make_model(
2     src_vocab, tgt_vocab, N=6, d_model=512, d_ff=2048, h=8, dropout=0.1
3 ):
4     "Helper: Construct a model from hyperparameters."
5     c = copy.deepcopy
6     attn = MultiHeadedAttention(h, d_model)
7     ff = PositionwiseFeedForward(d_model, d_ff, dropout)
8     position = PositionalEncoding(d_model, dropout)
9     model = EncoderDecoder(
10         Encoder(EncoderLayer(d_model, c(attn), c(ff), dropout), N),
11         Decoder(DecoderLayer(d_model, c(attn), c(attn), c(ff), dropout), N),
12         nn.Sequential(Embeddings(d_model, src_vocab), c(position)),
13         nn.Sequential(Embeddings(d_model, tgt_vocab), c(position)),
14         Generator(d_model, tgt_vocab),
15     )
16
17     # This was important from their code.
18     # Initialize parameters with Glorot / fan_avg.
19     for p in model.parameters():
20         if p.dim() > 1:
21             nn.init.xavier_uniform_(p)
22     return model

```