



$source\_padding\_idx \equiv source\_vocab\_size$   
 $target\_padding\_idx \equiv target\_vocab\_size$

These indices could be used to indicate the required padding when source and target language sentences are not of the same length.

$$\left. \begin{aligned} d_m \equiv d_{model} &= 200 \\ n \equiv num\_heads &= 4 \end{aligned} \right\} \Rightarrow d_n = \frac{d_m}{n} = 50$$

$$d_{ff} : d_{feed-forward} = 800$$

I) Source and Target Data Construction:

$source\_data.shape = torch.Size([20, 10])$   
 $target\_data.shape = torch.Size([20, 10])$

General shape of source and target data:

$$\begin{array}{ccc}
 [batch\_size \times max\_seq\_length] & & \\
 \parallel & & \parallel \\
 20 & & 10
 \end{array}$$

### II) Source and Target Masks Construction:

Transformer.py > generate\_mask(source, target)

```

source : [ batch-size x max-seq-length ]
target  : [ batch-size x max-seq-length ]

```

Breakdown computation of source and target masks:

```

SM = (source != self.source-padding-idx)
SM.shape = torch.Size([20, 10])
SM = SM.unsqueeze(1)
SM.shape = torch.Size([20, 1, 10])
SM = SM.unsqueeze(2)
SM.shape = torch.Size([20, 1, 1, 10])

```

```

TM = (target != self.target-padding-idx)
TM.shape = torch.Size([20, 10])
TM = TM.unsqueeze(1)
TM.shape = torch.Size([20, 1, 10])
TM = TM.unsqueeze(2)
TM.shape = torch.Size([20, 1, 1, 10])

```

SM : [ batch-size x 1 x 1 x max-seq-length ]

TM : [ batch-size x 1 x max-seq-length x max-seq-length ]

SL = target.size(1)

NPM = (1 - torch.triu(torch.ones(1, SL, SL, device = self.device), diagonal = 1)).bool()

NPM.shape = torch.Size([1, 10, 10])

↳ Elements above the main diagonal are False, while elements below the main diagonal are True.

TM = TM & NPM

TM.shape = torch.Size([20, 1, 10, 10])

(\*) When processing the source sequence all words are taken into consideration when extracting the MultiHeadAttention vector representation of each word. That is, each word can look at any word of the sequence.

(\*) When processing the target sequence each word can look only at the past instances of the sentence that is at the words that precede it.

### III) Encoder and Decoder Output Computation:

Transformer.py > forward(source, target)

source\_embedded.shape = torch.Size([20, 10, 200])

target\_embedded.shape = torch.Size([20, 10, 200])

General shape of embedded source and target sequences:

[	batch-size	x	max-seq-length	x	d-model	]
	20		10		200	

encoder\_output.shape = torch.Size([20, 10, 200])

Decoder output before the final fully connected layer:

decoder\_output.shape = torch.Size([20, 10, 200])

General shape of the encoder and decoder outputs:

[	batch-size	x	max-seq-length	x	d-model	]
	20		10		200	



General shape of the intermediate out tensor after the view operation:

`out : [batch-size * max-seq-length x target-vocab-size]`

`tgt = target-data.contiguous().view(-1)`

`tgt.shape = torch.Size([200])`

General shape of the intermediate tgt tensor after the view operation:

`tgt : [batch-size * max-seq-length]`

VI) Compute MultiHeadAttention Output within EncoderLayer and the first stage of the DecoderLayer { SelfAttention Mechanism }:

```
EncoderLayer.py > forward(x, mask) >
    attn_output = self.self_attn(x, x, x, mask)
where x ≡ source_embedded ≡ encoder_output
and mask ≡ source_mask
```

↓  
From previous encoder layers

```
DecoderLayer.py > forward(x, enc_output,
    source_mask, target_mask) >
    attn_output = self.self_attn(x, x, x, target_mask)
where x ≡ target_embedded ≡ decoder_output
```

↓  
From previous decoder layers

`EncoderLayer.py > forward(x, mask) >`

`attn_output = self.self_attn(x, x, x, mask)`

`x`  $\equiv$  source-embedded  $\equiv$  encoder-output

↓  
from previous  
encoder layers

`mask`  $\equiv$  source-mask

`x.shape = torch.Size([20, 10, 200])`

`x` : [ batch-size  $\times$  max-seq-length  $\times$  d-model ]

`mask.shape = torch.Size([20, 1, 1, 10])`

`mask` : [ batch-size  $\times$  1  $\times$  1  $\times$  max-seq-length ]

`MultiHeadAttention.py > forward(Q, K, V, mask)`

`W_q` : FULLY CONNECTED LINEAR LAYER [ d-model  $\times$  d-model ]

`W_k` : FULLY CONNECTED LINEAR LAYER [ d-model  $\times$  d-model ]

`W_v` : FULLY CONNECTED LINEAR LAYER [ d-model  $\times$  d-model ]

`W_o` : FULLY CONNECTED LINEAR LAYER [ d-model  $\times$  d-model ]

$$Q.\text{shape} = \text{torch.Size}([20, 10, 100])$$

$$K.\text{shape} = \text{torch.Size}([20, 10, 100])$$

$$V.\text{shape} = \text{torch.Size}([20, 10, 100])$$

$$Q = K = V = x \quad (\text{initial state})$$

$$Q, K, V : [\text{batch-size} \times \text{max-seq-length} \times d_{\text{model}}]$$

$$Q_0 = \text{self.W}_q(Q)$$

$$K_0 = \text{self.W}_k(K)$$

$$V_0 = \text{self.W}_v(V)$$

$$Q_0.\text{shape} = \text{torch.Size}([20, 10, 200])$$

$$K_0.\text{shape} = \text{torch.Size}([20, 10, 200])$$

$$V_0.\text{shape} = \text{torch.Size}([20, 10, 200])$$

$$Q_0, K_0, V_0 : [\text{batch-size} \times \text{max-seq-length} \times d_{\text{model}}]$$

$$Q \equiv x \longrightarrow \boxed{W_q : [200 \times 100]} \longrightarrow Q_0 : [20 \times 10 \times 200]$$

$$K \equiv x \longrightarrow \boxed{W_k : [200 \times 100]} \longrightarrow K_0 : [20 \times 10 \times 200]$$

$$V \equiv x \longrightarrow \boxed{W_v : [200 \times 100]} \longrightarrow V_0 : [20 \times 10 \times 200]$$

MultiHeadAttention.py > forward(Q, K, V, mask)  
 > split-heads(x)

$\left\{ \begin{array}{l} x \equiv Q_0 : [20 \times 10 \times 100] \\ x \equiv K_0 : [20 \times 10 \times 100] \\ x \equiv V_0 : [20 \times 10 \times 100] \end{array} \right.$

[ batch-size x max-seq-length x d-model ]

batch-size, seq-len, d-model = x.size()  
 ||                   ||                   ||  
 20                   10                   200

$Q_s = Q_0.view(\text{batch-size}, \text{seq-len}, \text{self.num-heads}, \text{self.d-k})$   
 $K_s = K_0.view(\text{batch-size}, \text{seq-len}, \text{self.num-heads}, \text{self.d-k})$   
 $V_s = V_0.view(\text{batch-size}, \text{seq-len}, \text{self.num-heads}, \text{self.d-v})$

$Q_s.shape = \text{torch.Size}([20, 10, 4, 50])$   
 $K_s.shape = \text{torch.Size}([20, 10, 4, 50])$   
 $V_s.shape = \text{torch.Size}([20, 10, 4, 50])$

$Q_s, K_s, V_s : [ \text{batch-size} \times \text{max-seq-length} \times \text{num-heads} \times \text{dk} ]$

```

Qt = Qs.transpose(1, 2)
Kt = Ks.transpose(1, 2)
Vt = Vs.transpose(1, 2)

```

```

Qt.shape = torch.Size([20, 4, 10, 50])
Kt.shape = torch.Size([20, 4, 10, 50])
Vt.shape = torch.Size([20, 4, 10, 50])

```

```

Qt, Kt, Vt: [batch-size x num-heads x max-seq-length x d.k]

```

```

MultiHeadAttention.py > forward(Q, K, V, mask)
Q ≡ Qt
K ≡ Kt
V ≡ Vt
> scaled-dot-product-attention(Q, K, V, mask)

```

```

Ktrans = Kt.transpose(-2, 1)

```

```

Ktrans.shape = torch.Size([20, 4, 50, 10])

```

General shape of Ktrans:

```

[batch-size x num-heads x d.k x max-seq-length]

```





Compute Attention Probabilities:

```
AP = torch.softmax(AT, -1)
```

```
AP.shape = torch.Size([20, 4, 10, 10])
```

General shape of AP at this stage:

[batch-size x num-heads x max-seq-len x max-seq-len]

⊗ We can verify that the sum of all elements across the last dimension of tensor AP are equal to 1 by invoking:  
torch.sum(AP, -1)

Compute Final Output for Scaled Dot Product Attention (within the context of the Encoder):

[20x4x10x10]    [20x4x10x50]  
                  ↑           ↑

```
A = torch.matmul(AP, Vt)
```

⊗ Tensor shapes are compatible in order for the multiplication to be executable!!

```
A.shape = torch.Size([20, 4, 10, 50])
```

General shape of A at this stage:

[batch-size x num-heads x max-seq-length x d-k]

```

MultiHeadAttention.py > forward(Q, K, V, mask)
                        > combine_heads(x)
output = self.W_o (self.combine_heads(attn_output))
x ≡ attn_output ≡ A

```

A.shape = torch.Size([20, 4, 10, 50])

0 1 2 3

General shape of x at this stage: {x ≡ A}

[batch-size × num-heads × max-seq-length × d-k]

```
Atrans = A.transpose(1, 2)
```

```
Atrans.shape = torch.Size([20, 10, 4, 50])
```

```
Atrans.is_contiguous == False
```

```
Atrans : [batch-size × max-seq-length × num-heads × d-k]
```

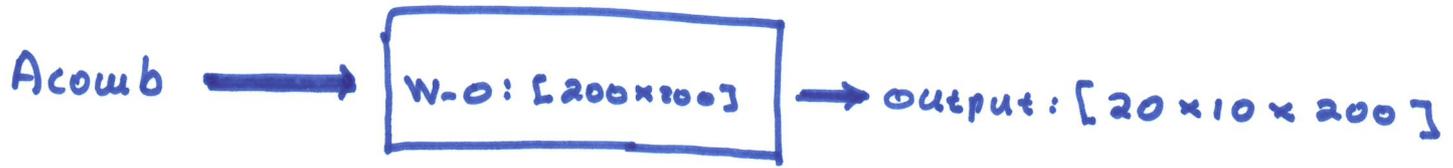
```
Acomb = Atrans.view(batch-size, seq-len, self.d_model)
```

```
Acomb.shape = torch.Size([20, 10, 200])
```

General shape of Acomb at this stage:

[batch-size × max-seq-length × d-model]

Compute the final output of the MultiHeadAttention Mechanism within the context of the Encoder:



General shape of output:

$[batch-size \times max-seq-length \times d-model]$

1718.  
Decoder.py > forward(x, enc-output, source-mask,  
target-mask)

attn-output = self.self\_attn(x, x, x, target-mask)  
x ≡ target-embedded ≡ decoder-output

x.shape = torch.Size([20, 10, 200])

x: [batch-size × max-seq-length × d-model]

target-mask.shape = torch.Size([20, 1, 10, 10])

target-mask: [batch-size × 1 × max-seq-length × max-seq-length]

Remaining steps are essentially the same as in  
the computation of the MultiHeadAttention Mechanism  
within the context of the Encoder.

The only difference is on the application of the  
target mask.

## Target Masking Application:

(19)

$AT = AT.$ masked\_fill (mask =  $\emptyset$ , -1e9)

mask  $\equiv$  target\_mask

$[20 \times 4 \times 10 \times 10]$

$[20 \times 1 \times 10 \times 10]$

General shape of AT at this stage:

$[batch\_size \times num\_heads \times max\_seq\_length \times max\_seq\_length]$

General shape of mask  $\equiv$  target\_mask at this stage:

$[batch\_size \times 1 \times max\_seq\_length \times max\_seq\_length]$

$AT.shape = torch.Size([20, 4, 10, 10])$

The previous operation is possible since the target\_mask tensor (mask of size:  $[20 \times 1 \times 10 \times 10]$ ) can be broadcast to a shape  $[20 \times 4 \times 10 \times 10]$  by expanding the dimension with size 1.

\* for  $i$  in range(batch\_size):

# The  $i$ -th slice of the target\_mask: target\_mask[i, :, :, :]

# will determine whether the  $i$ -th slice of the attn\_scores:

# attn\_scores[i, :, :, :] gets replaced with -1e9.

Compute Cross Attention:

```
Decoder.py > forward(x, enc_output, source_mask, target_mask)
```

```
attn_output = self.cross_attn(x, enc_output, enc_output, source_mask)
```

x ≡ target-embedded ≡ decoder-output

- x.shape = torch.Size([20, 10, 200])
- enc\_output.shape = torch.Size([20, 10, 200])
- source\_mask.shape = torch.Size([20, 1, 1, 10])
- target\_mask.shape = torch.Size([20, 1, 20, 10])
- x: [batch-size x max-seq-length x d-model]
- enc\_output: [batch-size x max-seq-length x d-model]
- source\_mask: [batch-size x 1 x 1 x max-seq-length]
- target\_mask: [batch-size x 1 x

```
MultiHeadAttention.py > forward(Q, K, V, mask)
```

```

Q ≡ x ≡ target-embedded ≡ decoder-output
K ≡ enc_output
V ≡ enc_output
mask ≡ source_mask

```

Based on the tensor assignments during calling the MultiHeadAttention mechanism in the context of the cross attention computation it is easy to deduce that:

$$\underline{\underline{S}} = \underline{\underline{Q}} \cdot \underline{\underline{K}}^T \quad (\text{per attention head})$$

↓
↓  
 decoder-output      encoder-output

$$\underline{\underline{\hat{S}}} = \text{softmax} \left( \sqrt{d_k}^{-1} \cdot \underline{\underline{S}} \right)$$

↓  
decoder-encoder outputs  
similarity values

$$\underline{\underline{Z}} = \underline{\underline{\hat{S}}} \cdot \underline{\underline{V}}$$

↓
↓
↓  
 encoder-output  
 decoder-encoder  
 similarities

↓

Provides an alternative representation of the source sequence taking into consideration its similarities with the target sequence.

## Positional Encoding Computation:

```
Transformer.py > forward(source, target)
```

```
source_embedded = self.dropout(self.positional_encoding(  
    self.encoder_embedding(source)))
```

```
target_embedded = self.dropout(self.positional_encoding(  
    self.decoder_embedding(target)))
```

```
source.shape = torch.Size([20, 10])  
target.shape = torch.Size([20, 10])
```

General shape of source and target tensors at this stage:  $[batch\_size \times max\_seq\_length]$ .

Acquire the embedding representation of the source and target tensors:

```
S = self.encoder_embedding(source)  
T = self.decoder_embedding(target)
```

```
S.shape = torch.Size([20, 10, 100])  
T.shape = torch.Size([20, 10, 100])
```

$S, T : [batch\_size \times max\_seq\_length \times d\_model]$

`PositionalEncoding.py > forward(x)`

$x \in S$  or  $x \in T$

```

batch_size, max_seq_length, d_model = x.size()
device = x.device

PE = torch.zeros(max_seq_length, d_model).to(device)
PE.shape = torch.Size([10, 200])

```

PE : [ max\_seq\_length x d\_model ]

```

position = torch.arange(0, max_seq_length, device=device,
                        dtype=torch.float)
position.shape = torch.Size([10])
position = position.unsqueeze(1)
position.shape = torch.Size([10, 1])

```

position : [ max\_seq\_length x 1 ]

```

numerator = torch.arange(0, d_model, 2, device=device,
                        dtype=torch.float)
numerator.shape = torch.Size([100])
numerator : [  $\frac{d\_model}{2}$  ]

```

↑ STEP

```

denominator = (torch.log(torch.FloatTensor([10000])).to(device) / d_model).to(device)
denominator.shape = torch.Size([1])

```

denominator : [ 1 ]

```
div-term = torch.exp(numerator * denominator)
div-term.shape = torch.Size([100])
```

div-term :  $\left[ \frac{d-model}{2} \right]$

Replace the even and odd positions of the positional encoding:

```
PE[:, 0::2] = torch.sin(position * div-term)
PE[:, 1::2] = torch.cos(position * div-term)
```

This operation could be further analyzed as:

```
even-positions = torch.sin([10x1] position * [100] div-term)
even-positions.shape = torch.Size([10, 100])
odd-positions = torch.cos([10x1] position * [100] div-term)
odd-positions.shape = torch.Size([10, 100])
```

even-positions, odd-positions :  $\left[ \text{max-seq-length} \times \frac{d-model}{2} \right]$

```
PE = PE.unsqueeze(0)
PE.shape = torch.Size([1, 10, 200])
PE : [ 1 x max-seq-length x d-model ]
```

```
PE0 = PE[:, :x.size(1)]
PE0.shape = torch.size([1, 10, 200])
PE0 : [ 1 x max-seq-length x d-model ]
```

```
X_pos = x + PE0
X_pos : [ batch-size x max-seq-length x d-model ]
X_pos.shape = torch.Size([20, 10, 200])
```