



# ★ Generate Random PyTorch Tensor:

## Transposition Operations

$Q = \text{torch.rand}(5, 4)$

$Q.$  shape

$\text{torch.Size}([5, 4])$

# ★ Perform Transposition Operations:

$Qa = Q.$  transpose  $(0, 1)$

$Qa.$  shape

$\text{torch.Size}([4, 5])$

$Qb = Q.$  transpose  $(-2, -1)$

$Qb.$  shape

$\text{torch.Size}([4, 5])$

$Qa-Qb = \text{tensor}([ [0, 0, 0, 0], [0, 0, 0, 0], [0, 0, 0, 0], [0, 0, 0, 0] ])$

## ★ Parameters:

batch-size = 4

seq-len = 8

d-model = 12

num-heads = 2

d-v = 6

max-seq-length = seq-len = 8

d-ff = 10

num-layers = 3

source-vocab-size = 20

target-vocab-size = 20

## Transposition and Viewing Operations

#13

batch-size = 4

seq-len = 8

d-model = 12

num-heads = 2

d-v = 6

$x = \text{torch.randn}(\text{batch-size}, \text{seq-len}, \text{d-model})$

$x.\text{shape}$

$\text{torch.Size}([4, 8, 12])$

$xv = x.\text{view}(\text{batch-size}, \text{seq-len}, \text{num-heads}, \text{d-v})$

$xv.\text{shape}$

$\text{torch.Size}([4, 8, 2, 6])$

$xt = xv.\text{transpose}(1, 2)$

$xt.\text{shape}([4, 2, 8, 6])$

Instantiate MultiHeadAttention module.

Input text:

batch-size = 4  
seq-len = 8  
d-model = 12  
num-heads = 2  
d.k = 6

from classes. MultiHeadAttention import MultiHeadAttention

mh = MultiHeadAttention(d\_model, num\_heads)

% set new Q tensor.

Q = torch.rand((batch\_size, seq\_len, d\_model))

K = torch.rand((batch\_size, seq\_len, d\_model))

V = torch.rand((batch\_size, seq\_len, d\_model))

% compute the attention output.

O = mh.scale\_dot\_product\_attention(Q, K, V)

% Make required imports

import torch, nn as nn

% Set the 4 Linear layers

QW = nn.Linear(d\_model, d\_model)

KW = nn.Linear(d\_model, d\_model)

VW = nn.Linear(d\_model, d\_model)

OW = nn.Linear(d\_model, d\_model)

% Pass tensors Q, K and V through the corresponding linear layers.

Q = QW(Q) → torch.Size([4, 8, 12])

K = KW(K) → torch.Size([4, 8, 12])

V = VW(V) → torch.Size([4, 8, 12])

Q.shape → torch.Size([4, 8, 12])  
K.shape → torch.Size([4, 8, 12])  
V.shape → torch.Size([4, 8, 12])  
O.shape → torch.Size([4, 8, 12])

% we could also call:

O = mh.forward(Q, K, V)

Linear(in\_features=12, out\_features=12, bias=True)

Checks the MultiHeadAttention

% Perform the head splitting operation

```

Qs = mha.split-heads(Q)
Ks = mha.split-heads(K)
Vs = mha.split-heads(V)

```

% Compute the scaled dot product output.

```

Os = mha.scaled-dot-product-attention(Qs, Ks, Vs) % Os.shape -> torch.Size([4, 2, 8, 6])

```

% Compute the combined version of O.

```

Oc = mha.combine-heads(Os) % Oc.shape -> torch.Size([4, 8, 12])

```

% Compute the final output of the attention layer.

```

Of = Ow(Oc) % Of.shape -> torch.Size([4, 8, 12])

```

★

### Checks Positional Encoding

from classes. PositionalEncoding import PositionalEncoding

# Set the maximum sequence length equal to the

# existing seq-len parameter.

max\_seq\_length = seq\_len

# Instantiate the PositionalEncoding class

PE = PositionalEncoding(d\_model, max\_seq\_length)

# Sample a random tensor in order to demonstrate

# the forward pass function call.

X = torch.rand((max\_seq\_length, d\_model)) # → x.shape = torch.Size([8, 12])

# Get the new version of X.

X = PE.forward(X) # → X.shape = torch.Size([1, 8, 12])

### Checks Positional Encoding

from classes. PositionalEncoding import PositionalEncoding

# Set the internal dimensionality parameter of the position-wise

# feed forward neural network.

d\_ff = 10

# Instantiate the PositionalEncoding class.

PWFF = PositionalEncoding(d\_model, d\_ff)

# Call the corresponding forward pass function.

X = PWFF.forward(X) # → X.shape = torch.Size([8, 12])

Mind that:

max\_seq\_length = 8 = seq\_len

d\_model = 12

d\_ff = 10

#E