

Input Transformation In Multi-Headed Attention.

Symbols:

d_{model} \rightarrow Input vector size (Embedding size for first layer)

d_k \rightarrow Length (size) of the Query, key and Value vectors in each head.

n_h \rightarrow Number of attention heads.

n_t \rightarrow Number of tokens in a sentence length.

t_{ij} \rightarrow Vector for token (j) in sentence (i)

q_{ijk} \rightarrow Query vector for token (j) in sentence (i) and head (k).

n_b \rightarrow batch size (Number of sentences in a batch)

k_{ijk} \rightarrow Key vector for token (j) in sentence (i) and head (k).

v_{ijk} \rightarrow Value vector for token (j) in sentence (i) and head (k).

a_{ijk} \rightarrow Attention score vector for token (j) in sentence (i) and head (k).
 \downarrow
 {from all other tokens in the sentence i.}

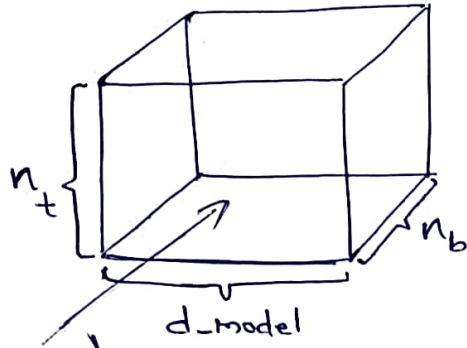
ah_{ijk} \rightarrow Attention head for token (j) in sentence (i) and head (k).

c_{ij} \rightarrow Concatenated Attention vector for token (j) in sentence (i).

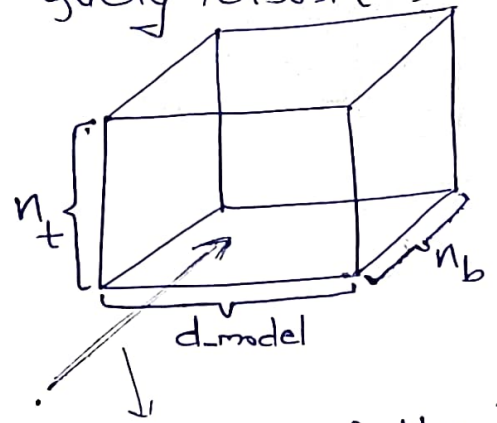
m_{ij} \rightarrow Output vector for token (j) in sentence (i) from the Multi-Headed attention layer.

Input Tensor (3D)

Query Tensor (3D)

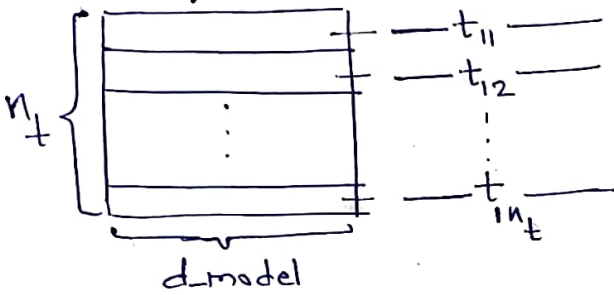


$nn.Linear(Input)$

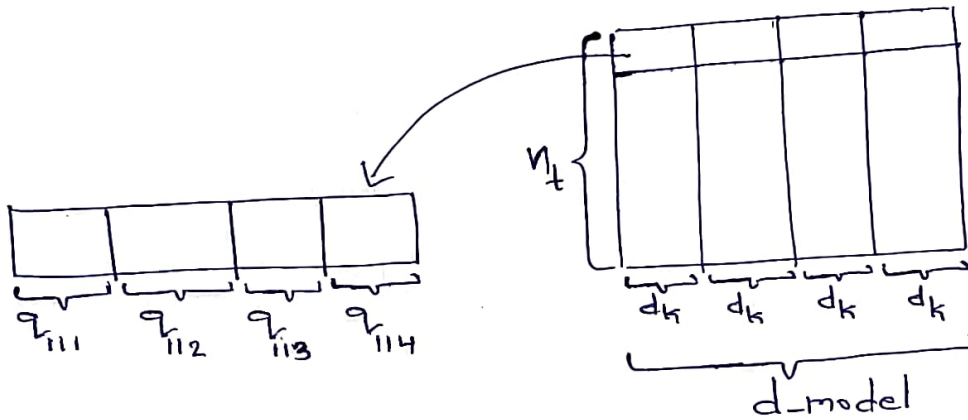
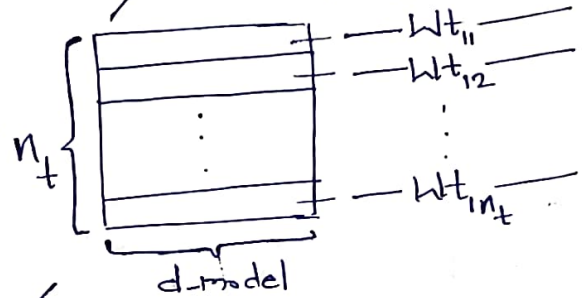


Front face of the Input tensor

Front face of the Query tensor



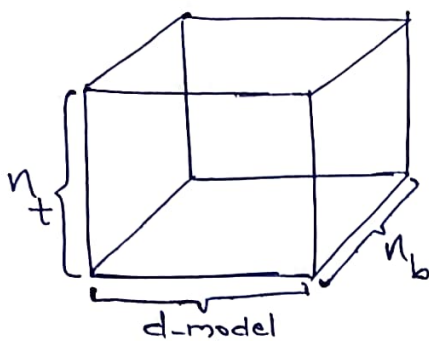
$nn.Linear(Input)$



Example case with 4 heads i.e., $n_h = 4$.

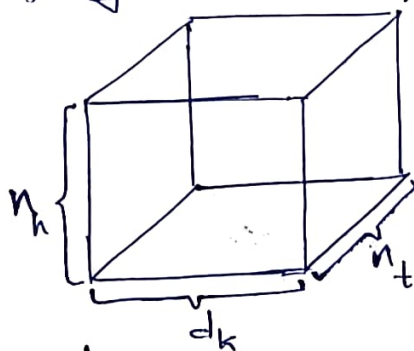
{ We are calculating query vectors for all the heads using a single matrix multiplication. }
(or a single $nn.Linear$ layer)

Query Tensor (3D)



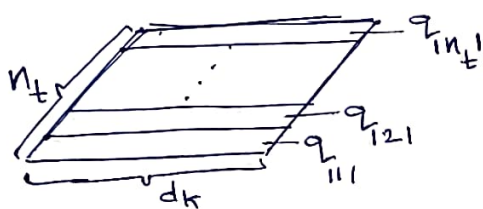
$tensor.view()$

Query Tensor Reshaped (4D)

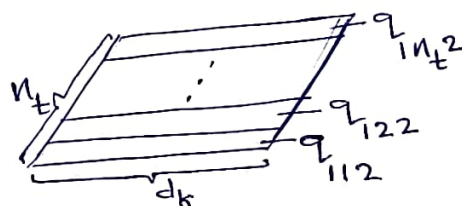


$n_b \rightarrow$ Number of these 3D tensors in the 4D Query tensor.

{ Omitted batch dimension since I cannot draw a 4D tensor }

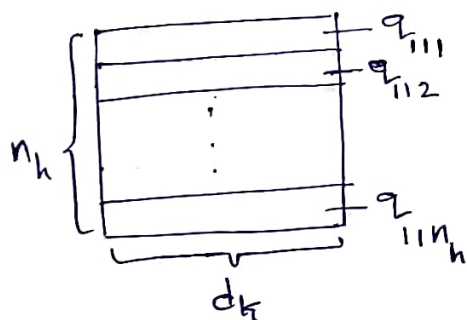


Top view of the first layer in the 4D Query tensor.

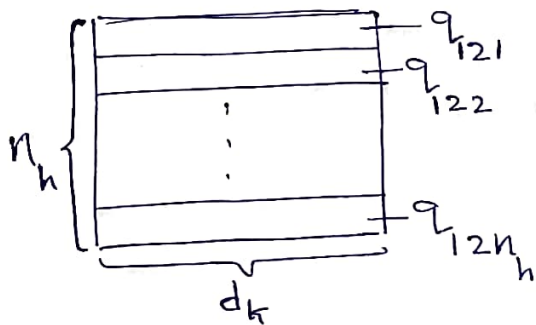


Top view of the second layer in the 4D Query tensor.

There will be n_h layers from Top view.



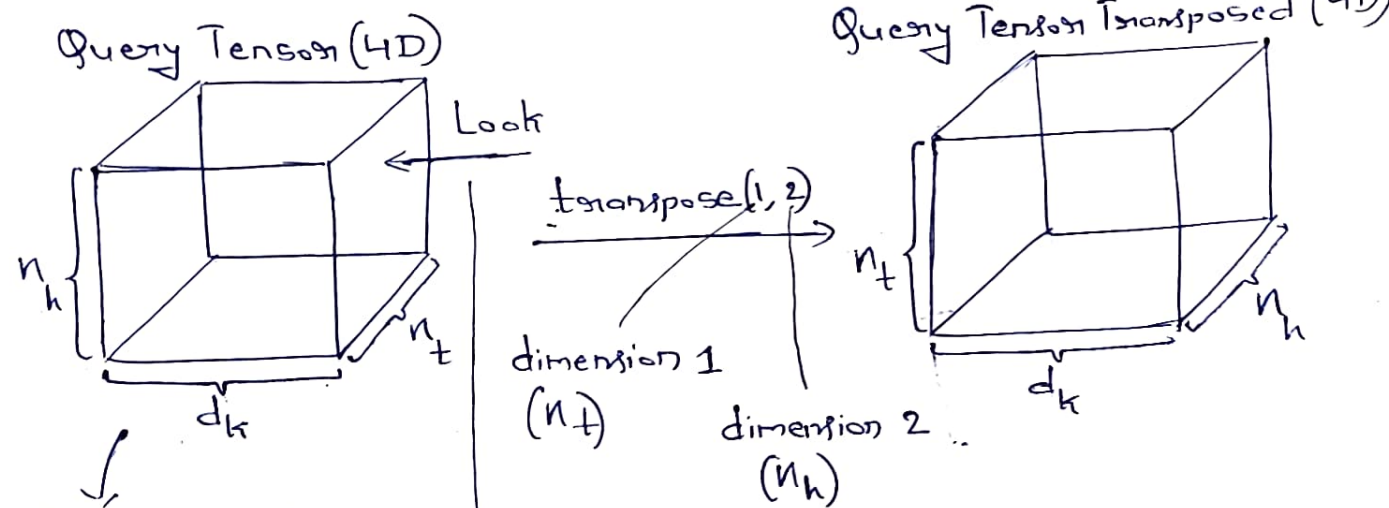
Front view of the first layer in the 4D Query tensor.



Front view of the second layer in the 4D Query tensor.

There are n_t layers from Front view.

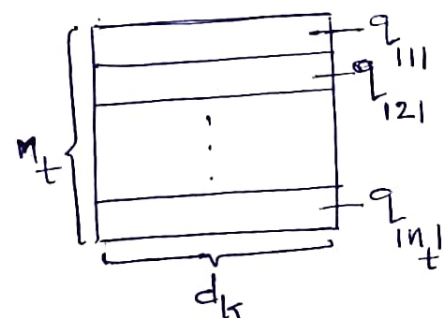
The same transformation is repeated for Key and Value tensors to obtain 4D tensors from 3D tensors.



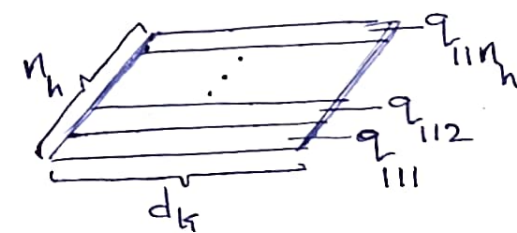
Omitted the batch dimension since I cannot draw a 4D tensor.

Look from the right side. This operation is equivalent to transposing the matrices we see when we look at the 3D tensors from right side.

The same transpose operation is again applied to the 4D key and value tensors.

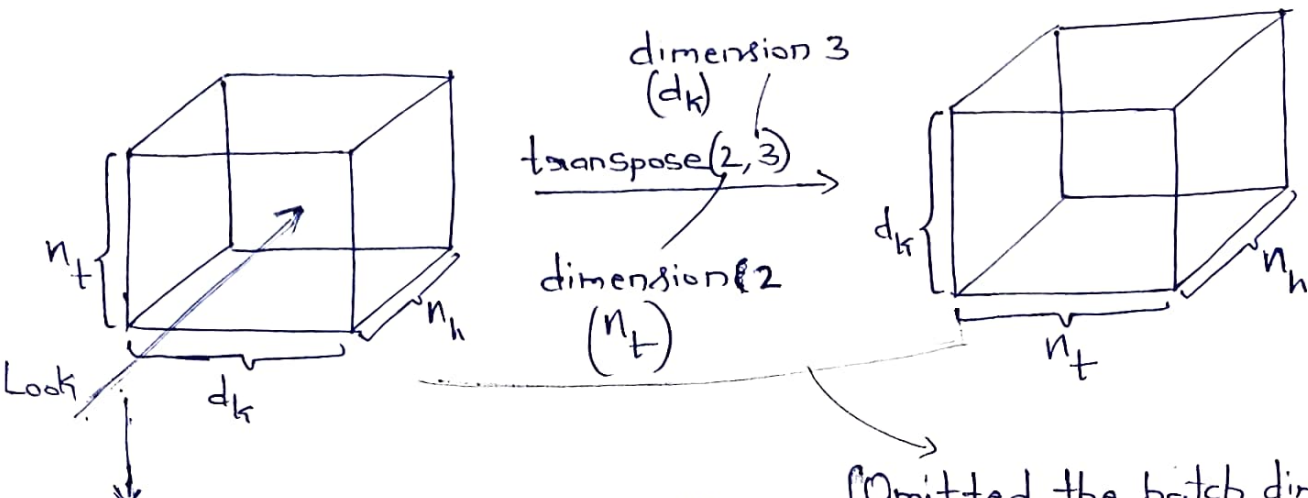


New front view of the Query after transpose operation



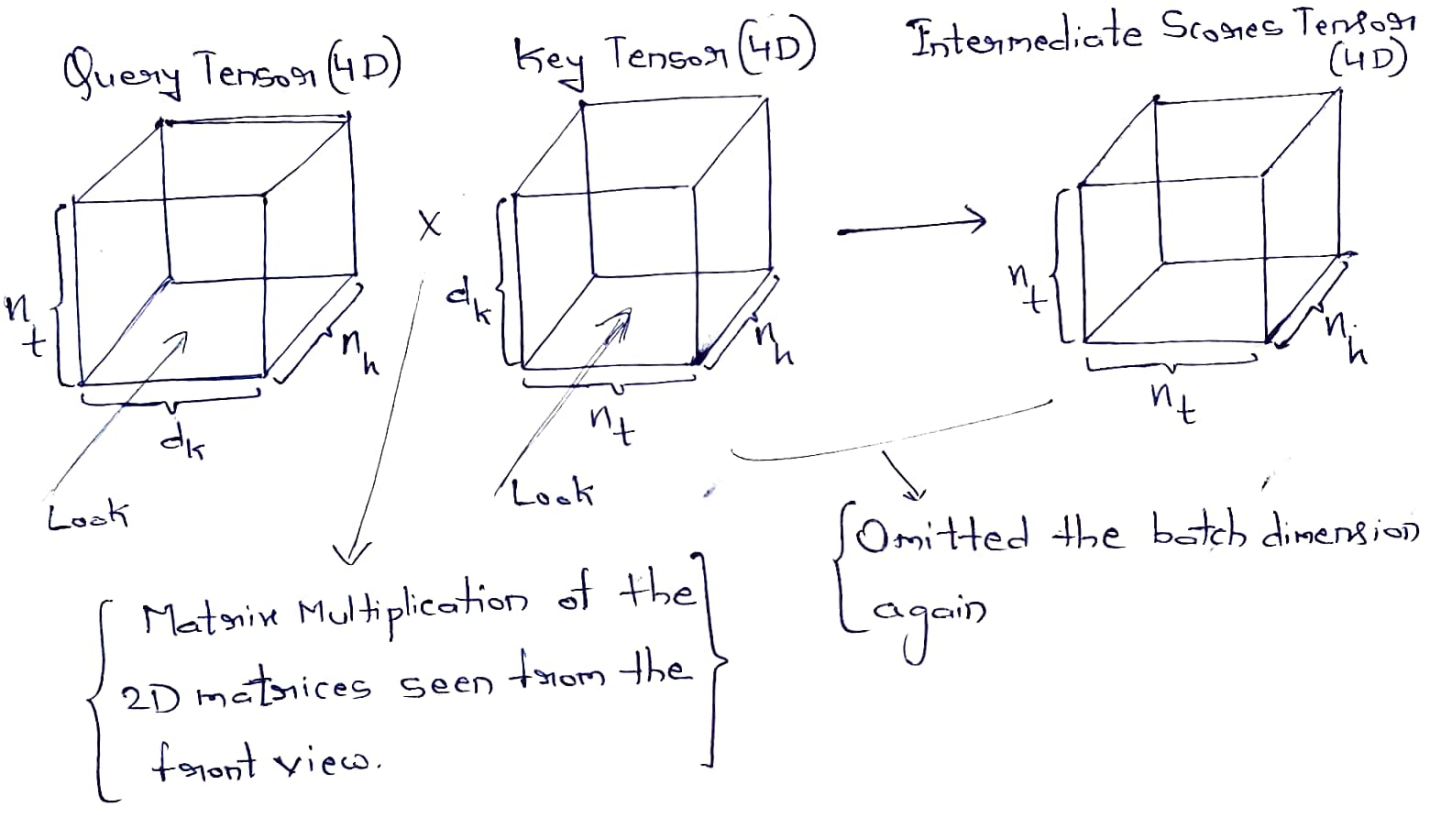
New front view of the Query after transpose operation.

Now, only the 4D key tensor is transposed again to calculate the attention scores.

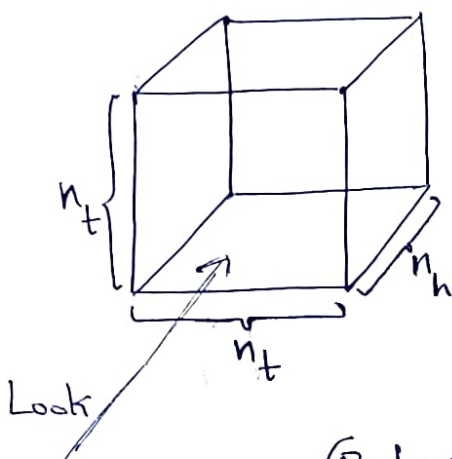


Look at the 3D tensor from front.
 This operation is equivalent to transposing the matrices we see when we look at the 3D tensor from front.

Omitted the batch dimension since I cannot draw a 4D tensor.

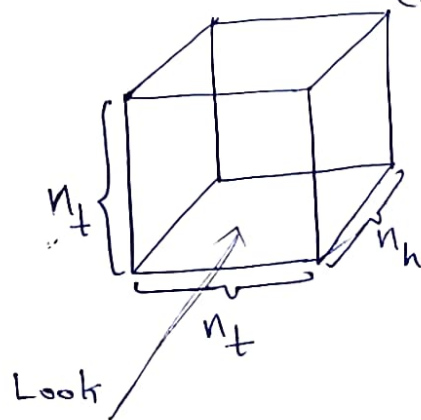


Intermediate Scores Tensor (4D)

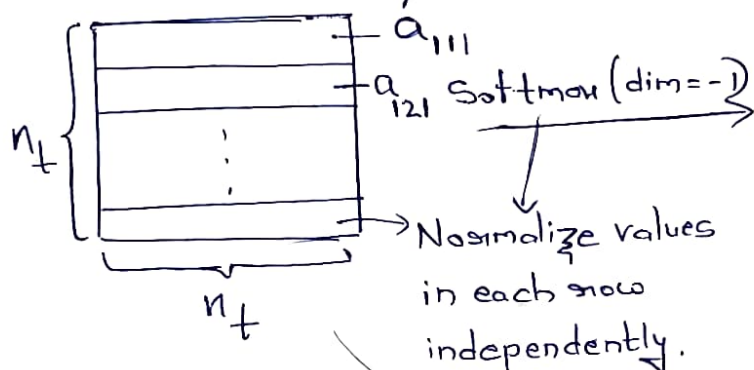


Softmax(dim=-1)
(n_t horizontal)

Attention Scores Tensor (4D) ⑥

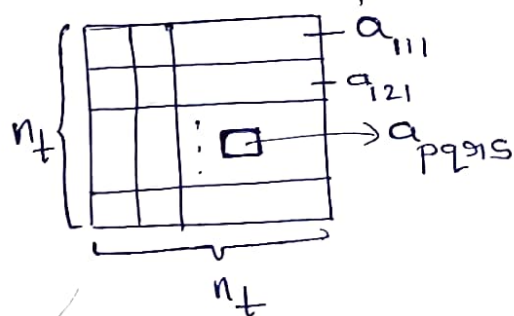


(Before Softmax)

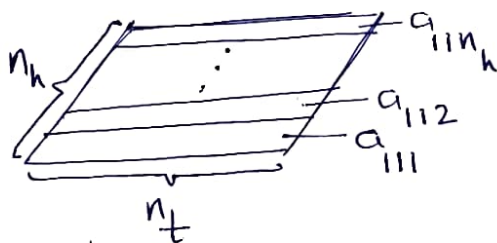


Normalize values
in each row
independently.

(Normalized scores)

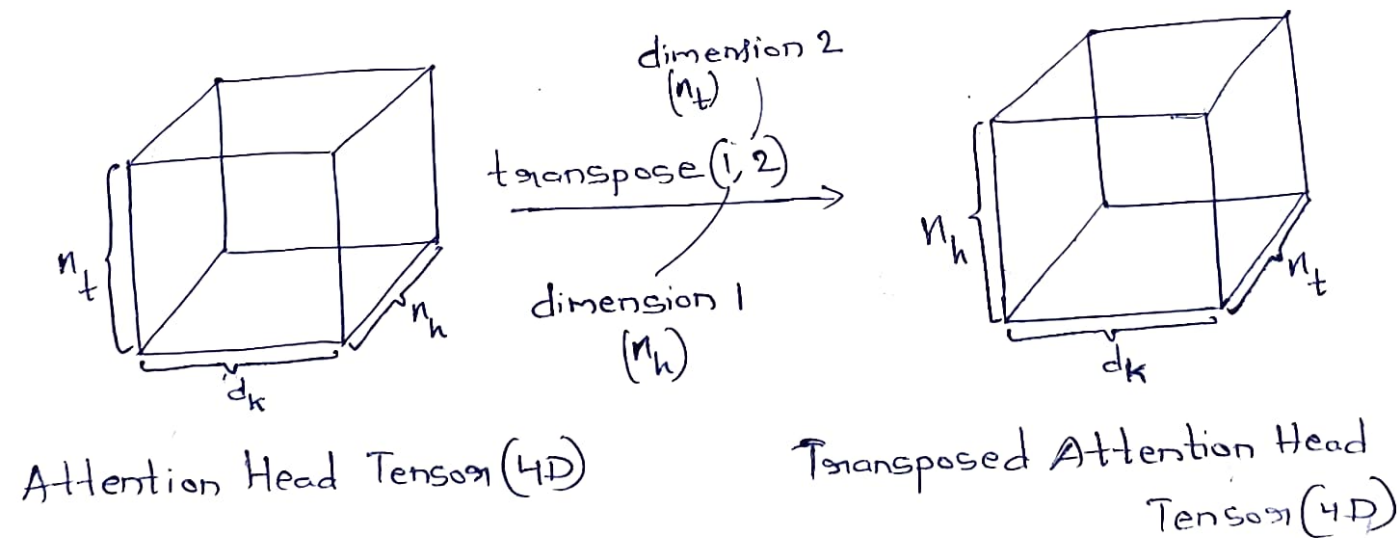
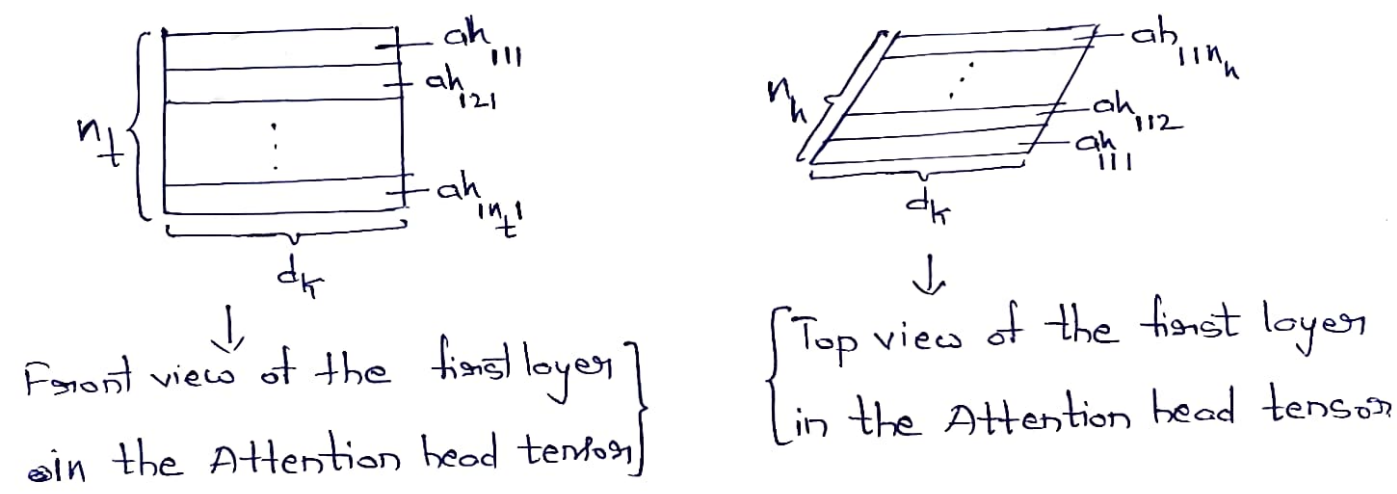
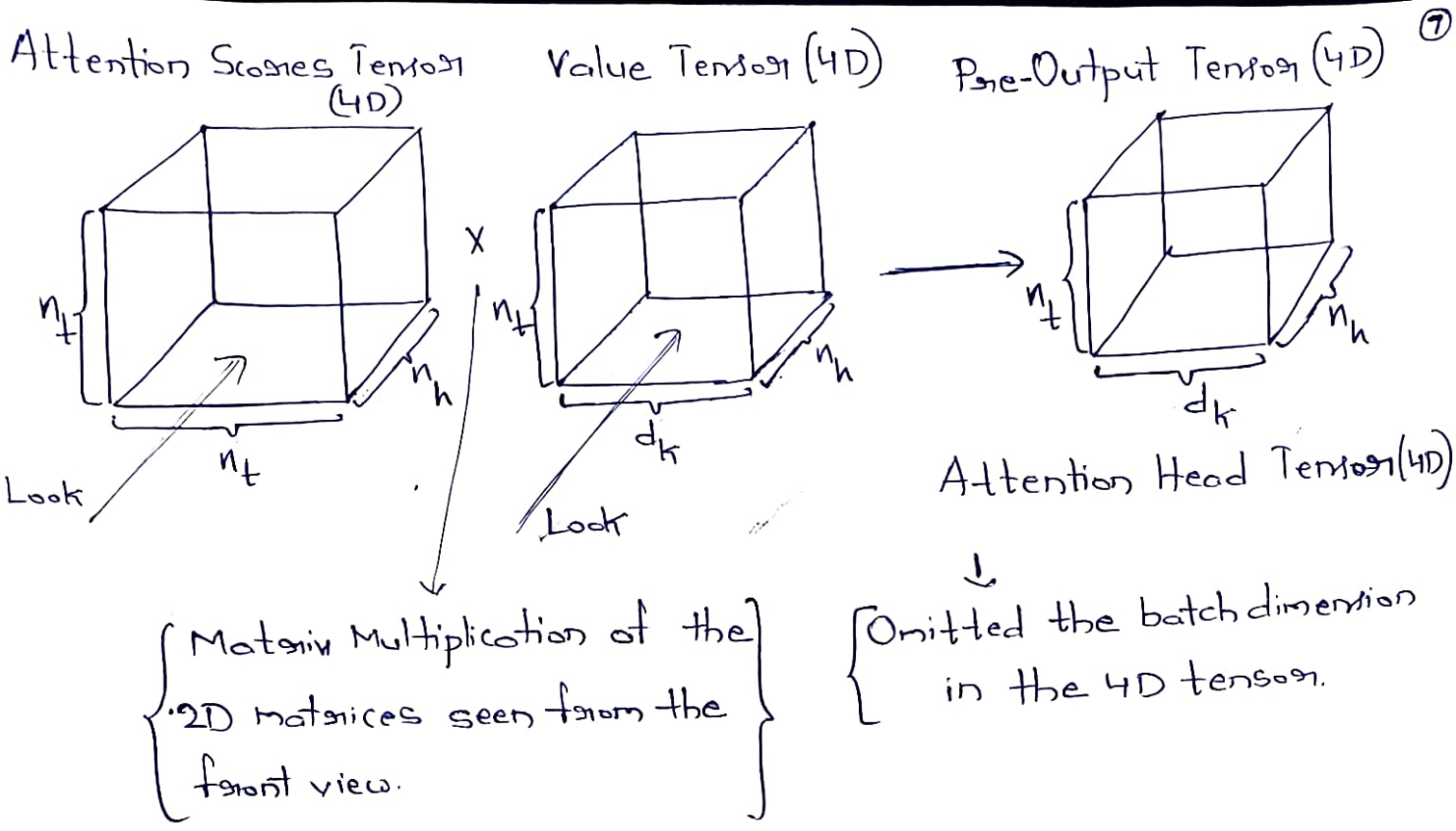


Front view of the first layer

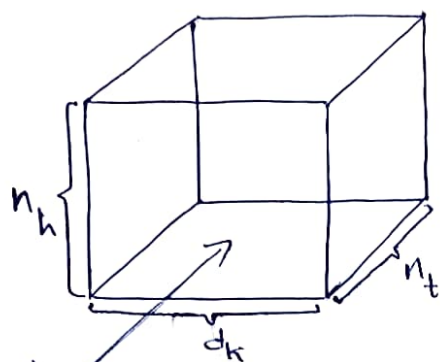


Top view of the first layer.

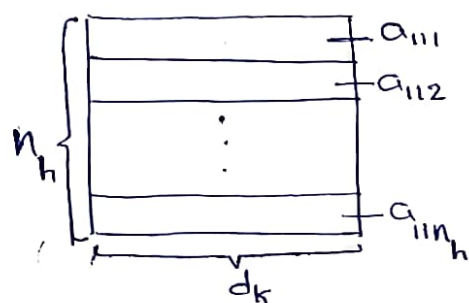
{Omitted the batch dimension for above}
4D tensors again.



Attention Head Tensor (4D)



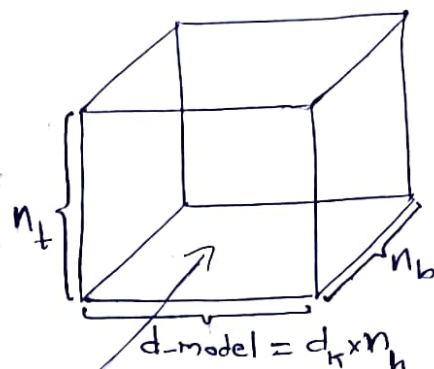
(Omitted, batch dimension)



Front view of the first layer of the Attention Head Tensor

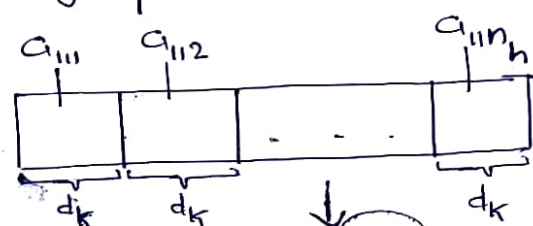
`tensor.view()`

Concatenated Attention Tensor (3D)

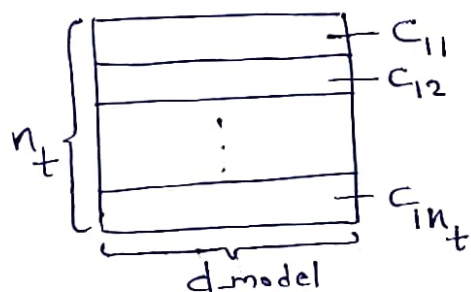


{ All attention heads for a single token are concatenated }

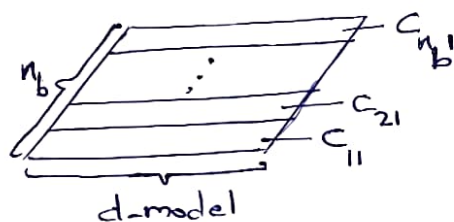
Concatenate all rows together



First row in the front view of the first layer of the Concatenated Attention Tensor

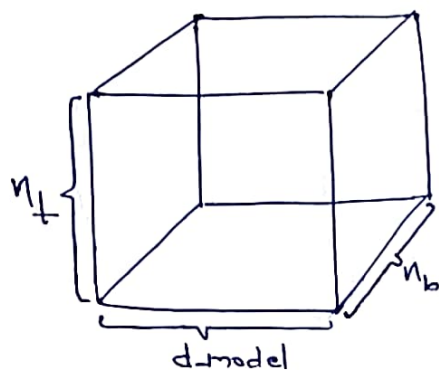


Front view of the first layer of the Concatenated Attention Tensor

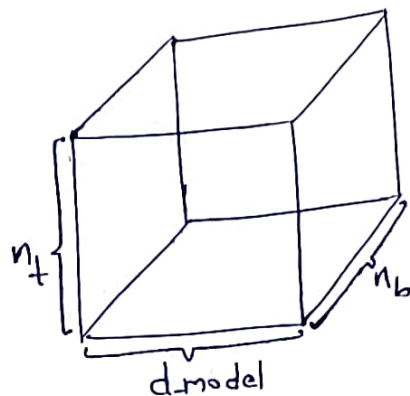


Top view of the first layer of the Concatenated Attention Tensor.

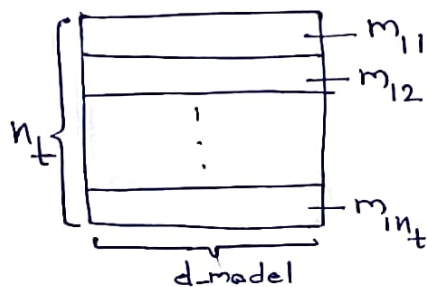
Concatenated Attention Tensor (3D)



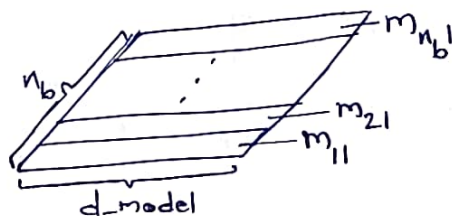
nn.Linear



Multi-Headed Attention
Layer Output Tensor (3D)



Front view of the first layer
of the Multi-Headed Attention
Layer Output Tensor



Top view of the first layer of
the Multi-Headed Attention
Layer Output Tensor.