

Consider an H -head attention, w/ k, q, v size d_1 , x size d .

How many parameters are introduced in 1 attention layer?

- ① Projection of x to k, q, v
- ② Compute score + softmax
- ③ Weighted avg of value vcs.
- ④ W^O projection

① W^q, W^k, W^v for 1 head
 $= [d \times d_1] - d d_1$

$3 \times H \times d d_1$

② 0

③ 0

~~④~~ ~~0~~

(4)

$W^0 \leftarrow$ input is

Concat of H v vectors

Input dim - $H \cdot d$

Output dim \rightarrow as is

same dim as $x_3 \rightarrow d$

$H d, d$

Overall $\rightarrow 3 H d, d + H d, d$

$= 4 H d, d$

Suppose (prev question)

+ L transformer layers

Vocabulary size V ,

FFNN hidden layer size

d_2 ,

Total no. of parameters??

→	Para.	in Attention Block	✓
→	"	" layer norm	
→	"	" FFNN	
→	"	"	⊔

$\times L$

QA \rightarrow

Input q "Who did...?"

Generate output tokens
from $M(\cdot | q)$

SA \rightarrow

Input $x \rightarrow$ Movie Review

Generate output tokens
from

$M(\cdot | \text{What is the sentiment
of } x?)$

Bi directional RNN

