

Vanilla attention: (original attention paper)

Assume input is : $(\underbrace{B}_{\text{batch size}}, \underbrace{T}_{\text{seq-len}}, \underbrace{D}_{\text{embedding}})$

We can ignore B for now and add this dimension at the end.

So (T, D) as input for the attention layer.





