

Pseudocode for Multi-Query Attention (MQA)

Siyu Yang

Algorithm 1 Uptraining Multi-Head to Multi-Query Attention

Require: M_{MHA} : Multi-Head Attention model

Require: D_{train} : Training dataset

Require: $\alpha = 0.05$: Proportion of original training compute for uptraining

```
1:  $M_{\text{MQA}} \leftarrow \text{ConvertToMQA}(M_{\text{MHA}})$ 
2: function CONVERTToMQA( $M_{\text{MHA}}$ )
3:   Initialize  $M_{\text{MQA}}$  with  $M_{\text{MHA}}$ 's architecture
4:   for each attention layer in  $M_{\text{MHA}}$  do
5:      $K_{\text{pooled}} \leftarrow \text{mean}(\{K_h | h \in \text{Heads}\})$ 
6:      $V_{\text{pooled}} \leftarrow \text{mean}(\{V_h | h \in \text{Heads}\})$ 
7:     Assign  $K_{\text{pooled}}, V_{\text{pooled}}$  to corresponding layer in  $M_{\text{MQA}}$ 
8:   end for
9:   return  $M_{\text{MQA}}$ 
10: end function
11:  $\text{steps}_{\text{total}} \leftarrow$  Number of steps in  $M_{\text{MHA}}$ 's original pre-training
12:  $\text{steps}_{\text{uptrain}} \leftarrow \lceil \text{steps}_{\text{total}} \times \alpha \rceil$ 
13: for  $\text{step} = 1$  to  $\text{steps}_{\text{uptrain}}$  do
14:    $\text{batch} \leftarrow$  Sample from  $D_{\text{train}}$ 
15:   Update  $M_{\text{MQA}}$  on  $\text{batch}$  ▷ Using pooled  $K$  and  $V$ 
16: end for
17: return  $M_{\text{MQA}}$ 
```
