

Pseudocode for Grouped-Query Attention (GQA) Uptraining

Siyu Yang, Based on the paper by Joshua Ainslie et al.

Algorithm 1 Uptraining for Grouped-Query Attention

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

Require: D_{train} : Training dataset

Require: G : Number of groups in GQA

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

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