

---

**Algorithm:** Fast data collection for MAP regression (interleaved pruning).

---

**Input:** pre-trained model  $Y_0$ , number of rounds  $rds$ , pruning budget  $k$ , total layers  $L$

```
1 Initialize train_data = {(\tilde{m}_{a,0}, \tilde{m}_{g,0}, a_0)};  
2 Function PruneIterative((\tilde{m}_{a,0}, \tilde{m}_{g,0}), (\tilde{m}_{a,max}, \tilde{m}_{g,max})):  
3   for n = 1 to rds do  
4     // Step 1: prune attention  
5      $\tilde{m}_{a,n} = \tilde{m}_{a,0} + n \cdot \frac{\tilde{m}_{a,max} - \tilde{m}_{a,0}}{rds};$   
6      $\tilde{m}_{g,n-1} = \tilde{m}_{g,0} + (n-1) \cdot \frac{\tilde{m}_{g,max} - \tilde{m}_{g,0}}{rds};$   
7     Prune  $Y_{n-1}$  along attention to ratio  $\tilde{m}_{a,n}$  to obtain  $Y'_n$ ;  
8     // Fine-tuning and evaluation  
9     Fine-tune  $Y'_n$  and evaluate  $\rightarrow (\tilde{m}_{a,n}, \tilde{m}_{g,n-1}, a_n)$ ;  
10    Append  $(\tilde{m}_{a,n}, \tilde{m}_{g,n-1}, a_n)$  to train_data;  
11    // Step 2: prune activation  
12     $\tilde{m}_{g,n} = \tilde{m}_{g,0} + n \cdot \frac{\tilde{m}_{g,max} - \tilde{m}_{g,0}}{rds};$   
13    Prune  $Y'_n$  along activation to ratio  $\tilde{m}_{g,n}$  to obtain  $Y_n$ ;  
14    // Fine-tuning and evaluation  
15    Fine-tune  $Y_n$  and evaluate  $\rightarrow (\tilde{m}_{a,n}, \tilde{m}_{g,n}, a_n)$ ;  
16    Append  $(\tilde{m}_{a,n}, \tilde{m}_{g,n}, a_n)$  to train_data;  
17  end  
18 Set  $\tilde{m}_{a,max} = \frac{k}{L}$  and  $\tilde{m}_{g,max} = \frac{k}{L}$ ;  
19 PruneIterative((\tilde{m}_{a,0}, \tilde{m}_{g,0}), (\tilde{m}_{a,max}, \tilde{m}_{g,max}));  
20 return train_data
```

---