**Algorithm 1** The pipeline of PTQ1.61

**Require:** Pretrained LLM $M$, input activation $X$
**Ensure:** Quantized LLM $M_q$

1: $M_P = \text{P}(M)$ ▷ Preprocessing
2: **for** each block $b$ in $M_P$ **do**
3:      $m_s \leftarrow \text{Top}_{20\%}(\text{mean}(|X[:, c, \ldots]|))$
4:      **for** each channel $c$ in $W$ **do**
5:          **if** $m_s[c] = 1$ **then**
6:              Quantize $(W_c)$ ▷ Salient weights
7:          **else**
8:              Binarize $(W_c)$ ▷ Unsalient weights
9:          **end if**
10:      **end for**
11:      $Y_{\text{full}} \leftarrow b(X)$
12:      $Y_q \leftarrow b_q(X)$
13:      $Y_{fq} \leftarrow b(Y_{q-1})$
14:      $\text{L}_{norm} = \text{MSE}(Y_q, Y_{full}) + \text{NLC}(Y_q, Y_{full})$
15:      $\text{L}_{homo} = \text{MSE}(Y_q, Y_{fq}) + \text{NLC}(Y_q, Y_{fq})$
16:      $\text{Loss} = \text{L}_{norm} + \text{L}_{homo}$
17:      Loss.backward() ▷ Optimize scaling factors of quantized weights
18: **end for**
19: **return** Quantized model $M_q$