**Algorithm 1** Ghost Backdoor based on Neuron Select

---

**Input:** central server $C_s$, a set of all client $C$, end epoch $E_e$, current client $C_i$, learning rate $\eta$, dataset $D$, mask matrix $R_{mask}^{r \times d}$, ghost neurons' values matrix $R_{V_s}^{r \times d}$

**Output:** a global model with high accuracy, stealth backdoor and high accuracy in main-task

---

1: $C_s$ select $n$ clients by random into $C_m$
2: $C_s$ build a global model $G$
3: $C_s$ send $G$ to each client in $C_m$
4: choose the ghost neurons
5: pre-train with benign samples to collect the values of every neurons
6: choose $V_s$ as trigger
7: **for** epoch $< E_e$ **do**
8:     **for** number $k$ of client in $C_m$ **do**
9:         Download $G$ as local model $L$ and train $L$ by benign datasets $D$,
10:         Compute gradient by $D$ on batch $B_i$ of size $\ell$
11:         $g_i = \frac{1}{n} \sum_{i=1}^{n} \nabla_\theta \mathcal{L}(\theta_{C_i}, D)$
12:         **if** client $C_i$ is advisary **and** epoch mod $N_{attack} = 0$ **then**
13:             $\hat{g}_i = g_i * R_{mask}^{r \times d} + R_{V_s}^{r \times d}$
14:             Update $\theta_{C_{i+1}} = \theta_{C_i} - \eta \hat{g}_i$
15:         **else**
16:             Update $\theta_{C_{i+1}} = \theta_{C_i} - \eta g_i$
17:         Upload $\theta_{C_{i+1}}$ to $C_s$
18:         $C_s$ recieve $\sum_1^k \theta_{C_{i+k}}$ and generate update gradient $U$ for $G$
19: $G_{i+1} = G_i - U_i$
20: **return** Final global model $G$ with backdoor