**Algorithm 1** Ghost Backdoor based on neuron select

---

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

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

1: $C_s$ select $n$ clients by random into $C_n$
2: $C_s$ build a global model $G$
3: $C_s$ send $G$ to each client in $C_n$
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** $e < E_e$ **do**
8:     **for** the $k$-th client $C_e^k$ in $C_n$ **do**
9:         Download $G$ as local model $L$ and train by $D$,
10:         Compute gradient by $D$ on batch $B_i$ of size $\ell$
11:         $g_{e+1}^k = \frac{1}{\ell} \sum_{i=1}^{\ell} \nabla_\theta \mathcal{L}(\theta_{C_e^k}, D)$
12:         **if** client $C_i$ is advisary **and** epoch mod $N_{attack}$ = 0 **then**
13:             $\hat{g}_{e+1}^k = g_{e+1}^k * \mathbb{R}_{mask}^{r \times d} + \mathbb{R}_{V_s}^{r \times d}$
14:             Update $\theta_{C_{e+1}^k} = \theta_{C_e^k} - \eta \hat{g}_{e+1}^k$
15:         **else**
16:             Update $\theta_{C_{e+1}^k} = \theta_{C_e^k} - \eta g_{e+1}^k$
17:         Upload $\theta_{C_{e+1}^k}$ to $C_s$
18:     $C_s$ recieve $\sum_1^k \theta_{C_{e+1}^k}$ and generate update $U_{e+1}$ for $G_{last}$
19:     $G = G_{last} - U_{e+1}$
    **return** Final global model $G$ with backdoor

---