

---

**Algorithm 1** Stealthing and Robust Backdoor based on Steganographic Algorithm

---

**Input:** start epoch  $E_s$ , attack num  $E_a$ , end epoch  $E_e$ , client set  $C$ , selected client set  $C_m$ , adversary set  $C_{adv}$ , global model  $G$ , local model  $\theta$ , central server  $C_s$ , benign datasets  $\hat{D}$ , poisoned datasets  $\hat{D}_p$ , benign learning rate  $\eta_b$ , poison learning rate  $\eta_p$ , Sparse-update gradient removal scale  $\mathcal{R}\%$

**Output:** a global model with high accuracy, stealth and robust 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: **for** epoch  $< E_e$  and epoch  $> E_s + E_a$  **do**
  - 5:     **for** number  $k$  of client in  $C_m$  **do**
  - 6:         **if** client  $e_i \in C_{adv}$  **then**
  - 7:             Download  $G$  as local model  $L$  and train  $L$  by poisoned datasets  $\hat{D}_p$ ,
  - 8:             Compute gradient by  $\hat{D}_p$  on batch  $B_i$  of size  $\ell$
  - 9:              $g_i^p = \frac{1}{\ell} \sum_{i=1}^{\ell} \nabla_{\theta} \mathcal{L}(\theta_{e_i}, \hat{D}_p)$
  - 10:              $\theta_{e_{i+1}} = \theta_{e_i} - \eta_p g_i^p$  where  $top5\%(Value(g)) \not\subseteq g_i^p$
  - 11:             set  $\mathcal{R}\%$  of gradient  $\theta_{e_{i+1}}$  to zero
  - 12:             Upload  $\theta_{e_{i+1}}$  to  $C_s$
  - 13:         **else if** client  $e_i \notin C_{adv}$  **then**
  - 14:             Download  $G$  as local model  $L$  and train  $L$  by private benign dataset  $\hat{D}$ ,
  - 15:             Compute gradient by  $\hat{D}_b$  on batch  $B_i$  of size  $\ell$
  - 16:              $g_i^b = \frac{1}{\ell} \sum_{i=1}^{\ell} \nabla_{\theta} \mathcal{L}(\theta_{e_i}, \hat{D})$
  - 17:              $\theta_{e_{i+1}} = \theta_{e_i} - \eta_b g_i^b$
  - 18:             Upload  $\theta_{e_{i+1}}$  to  $C_s$
  - 19:      $C_s$  receive  $\sum_1^k \theta_{e_{i+1}}^k$  and generate update gradient  $U$  for  $G$
  - 20:      $G_{i+1} = G_i - U_i$
  - 21: **return** Final global model  $G$  with backdoor
-