---

**Algorithm 1: Server Aggregation**

---

**Inputs:** $\mathbf{w}'_g$: the global model of last FL training round; $\mathcal{W}_l$: the list of local models submitted by each client in the current FL training round.

**Variables:** $\mathcal{A}$: A FedAttacker instance initialized based on the FL configuration file; $\mathcal{D}$: A FedDefender instance that is initialized based on the FL configuration file.

1 **Function** $server\_aggregation(\mathcal{W}_l)$ **begin**
2     $\mathcal{W}_l \leftarrow before\_aggregation\_process(\mathcal{W}_l, \mathbf{w}'_g)$
3     $\mathbf{w}_g \leftarrow before\_aggregation\_process(\mathcal{W}_l, \mathbf{w}'_g)$
4     **return** $after\_aggregation\_process(\mathcal{W}_l, \mathbf{w}_g)$

5 **Function** $before\_aggregation\_process(\mathcal{W}_l, \mathbf{w}'_g)$ **begin**
6     **if** $\mathcal{A}.is\_attack\_enabled()$ **then**
7        **if** $\mathcal{A}.is\_data\_reconstruction\_attack()$ **then** $\mathcal{A}.reconstruct\_data(\mathcal{W}_l, \mathbf{w}'_g)$ ;
       **if** $\mathcal{A}.is\_model\_poisoning\_attack()$ **then** $\mathcal{W}_l \leftarrow \mathcal{A}.poison\_model(\mathcal{W}_l, \mathbf{w}'_g)$;
8     **if** $\mathcal{D}.is\_defense\_enabled()$ & $\mathcal{D}.is\_defense\_before\_aggregation()$ **then**
       $\mathcal{W}_l \leftarrow \mathcal{D}.defend\_before\_aggregation(\mathcal{W}_l, \mathbf{w}'_g)$
9     **return** $\mathcal{W}_i$

10 **Function** $on\_aggregation\_process(\mathcal{W}_l, \mathbf{w}_g)$ **begin**
11     **if** $\mathcal{D}.is\_defense\_enabled()$ & $\mathcal{D}.is\_defense\_on\_aggregation()$ **then**
       **return** $\mathcal{D}.defend\_on\_aggregation(\mathcal{W}_l, \mathbf{w}_g)$
12     **return** $aggregate(\mathcal{W}_i)$

13 **Function** $after\_aggregation\_process(\mathbf{w}_g)$ **begin**
14     **if** $\mathcal{D}.is\_defense\_enabled()$ & $\mathcal{D}.is\_defense\_after\_aggregation()$ **then**
       **return** $\mathcal{D}.defend\_after\_aggregation(\mathbf{w}_g)$
15     **return** $\mathbf{w}_g$

---

---

**Algorithm 2: Client Training**

---

**Inputs:** $dataset$: the local dataset of a client.

**Variables:** $\mathcal{A}$: A FedAttacker instance initialized based on the FL configuration file;

1 **Function** $client\_training(dataset)$ **begin**
2     **if** $\mathcal{A}.is\_attack\_enabled()$ & $\mathcal{A}.is\_data\_poisoning\_attack()$ **then**
       $dataset \leftarrow \mathcal{A}.poison\_data(dataset)$
3     $\mathbf{w}_l \leftarrow train(dataset)$
4     $send\_to\_server(\mathbf{w}_l)$

---