

DeepSight: Mitigating Backdoor Attacks in Federated Learning Through Deep Model Inspection

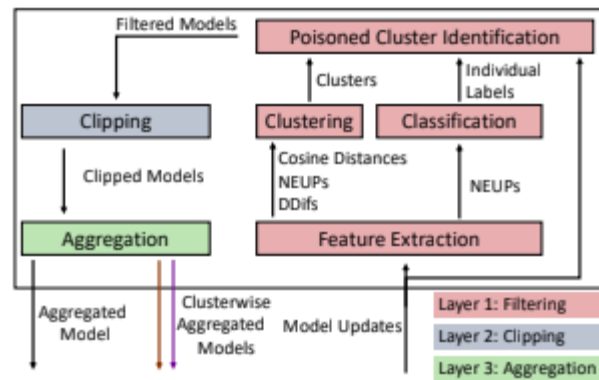
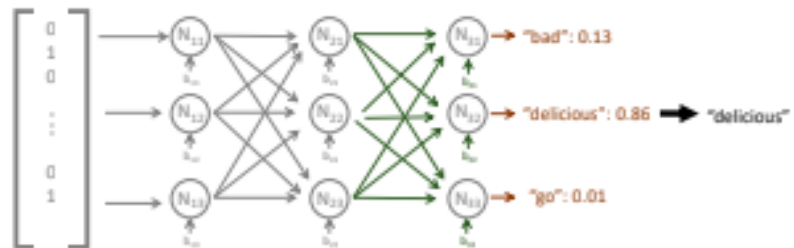


Fig. 2: Structure of DeepSight

Algorithm 1 Filtering Layer

```
1: Input:  $N$ , ▷ number of models
2:  $W_i$  ▷ list of  $N$  received local models
3:  $G_t$  ▷ global model
4: Parameters:  $\tau$ , ▷ Threshold of suspicious models for excluding cluster
5: seeds, ▷ 3 seeds for generating random data for ddifs
6: input_dim ▷ dimension of a single input
7: Output: accepted_models
8: ▷ Feature Extraction
9: cosine_distances  $\leftarrow 0^{N \times N}$ 
10: global_bias  $\leftarrow$  output_layer_bias( $G_t$ )
11: for each clients  $i, j$  in  $[1, N]$  do
12:   update $_i$   $\leftarrow$  output_layer_bias( $W_i$ ) - global_bias
13:   update $_j$   $\leftarrow$  output_layer_bias( $W_j$ ) - global_bias
14:   cosine_distances $_{i,j}$   $\leftarrow 1 - \text{COSINE}(\text{update}_i, \text{update}_j)$ 
15: end for
16:  $\forall i \in \{1, \dots, N\}$  : neups $_i$   $\leftarrow$  NEUPs( $G_t, W_i$ )
17:  $\forall i \in \{1, \dots, N\}$  : thresh_exds $_i$   $\leftarrow$  THRESHOLD_EXCEEDING(neups)
18:  $\forall i \in \{1, 2, 3\}$  : rand_input_data $_i$   $\leftarrow$  random_matrix(seeds $_i$ , 20000,
    input_dim)
19:  $\forall i \in \{1, 2, 3\}$  : ddifs $_i$   $\leftarrow$  DDIFs(rand_input_data $_i, G_t, W_1 \dots W_n$ )
20: ▷ Classification
21: classificat_boundary  $\leftarrow$  MEDIAN(thresh_exds) / 2
22:  $\forall i \in \{1, \dots, N\}$  : labels $_i$   $\leftarrow$  (thresh_exds $[i] \leq$  classificat_boundary)? 1:0
23: ▷ Clustering
24: clusters  $\leftarrow$  CLUSTER( $N$ , neups, ddifs, cosine_distances)
25: ▷ PCI
26: accepted_models  $\leftarrow \{\}$ 
27: for cluster in clusters do
28:   amount_of_positives  $\leftarrow$  SUM(labels[cluster]) / |cluster|
29:   if amount_of_positives  $< \tau$  then
30:     accepted_models  $\leftarrow$  accepted_models  $\cup$  models[cluster]
31:   end if
32: end for
```

Algorithm 2 Clustering

```
1: procedure DISTSFROMCLUST(clusters, N)
2:    $\forall i, j \in \{1, \dots, N\} : \text{pairwise\_dists}_{i,j} \leftarrow \text{cluster\_of\_model}(i, \text{clusters})$ 
    $\text{== cluster\_of\_model}(j, \text{clusters})? 0:1$   $\triangleright \text{cluster\_of\_model}(x, \text{clusters})$ 
   returns the cluster that contains the model with index x
3:   return pairwise_dists
4: end procedure
5:
6: Input:
7: N,  $\triangleright N$  is the number of models
8: neups,  $\triangleright$  NEUPs as list of N vectors with dimension P
9: ddifs  $\triangleright$  DDifs as list of 3 lists of vectors with dimension P
10: cosine_distances  $\triangleright$  cosine_distances a matrix  $\in \mathbb{R}^{N \times N}$ 
11: Output: clusters  $\triangleright$  clusters as set of sets of indices
12:
13: cosine_clusters  $\leftarrow$  HDBSCAN(distances = cosine_distances)
14: cosine_cluster_dists  $\leftarrow$  DistsFromClust(cosine_clusters, N)
15: neup_clusters  $\leftarrow$  HDBSCAN(values = neups)
16: neup_cluster_dists  $\leftarrow$  DistsFromClust(NEUP_clusters, N)
17:  $\forall i \in \{1, 2, 3\} : \text{ddif\_clusters}_i \leftarrow$  HDBSCAN(values = ddifsi)
18:  $\forall i \in \{1, 2, 3\} : \text{ddif\_clust\_dists}_i \leftarrow$  DistsFromClust(ddif_clustersi, N)
19: merged_ddif_clust_dists  $\leftarrow$  AVG(ddif_clust_dists1, ddif_clust_dists2,
   ddif_clust_dists3)
20:  $\triangleright$  Combine clusterings
21: merged_distances  $\leftarrow$  AVG(merged_ddif_clust_dists, neup_clust_dists, co-
   sine_clust_dists)
22: clusters  $\leftarrow$  HDBSCAN(distances = merged_distances)
```
