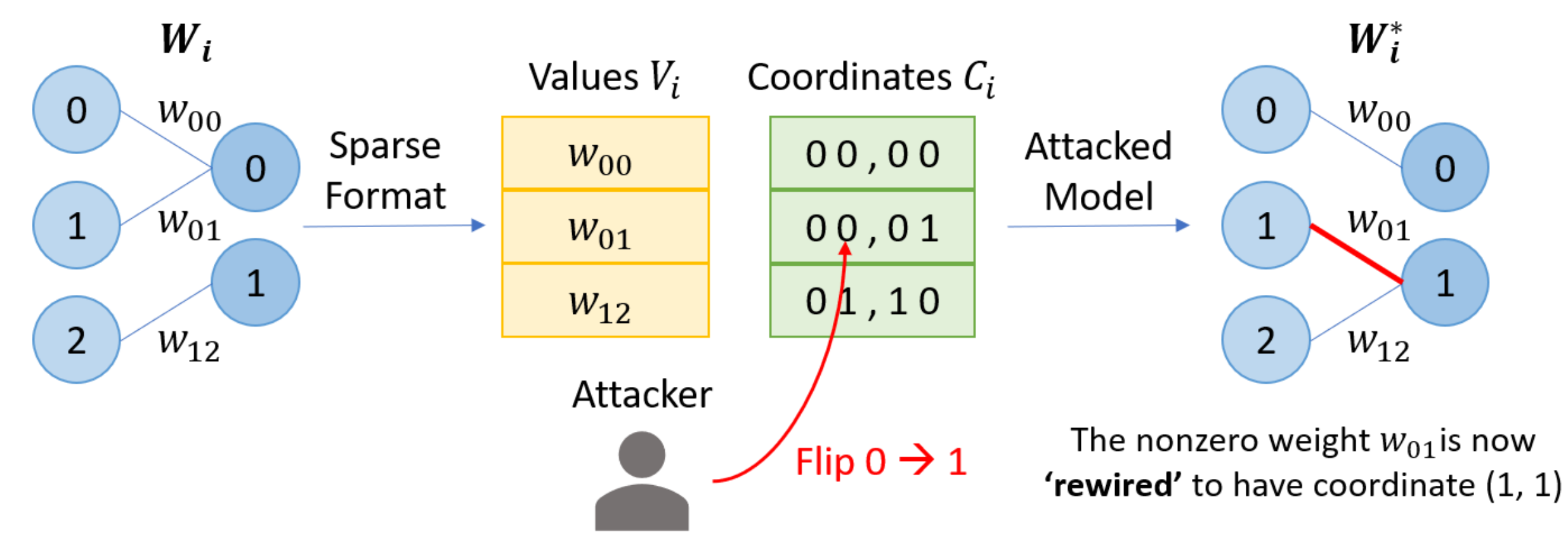


Introduction

- Exploiting **hardware-level faults** to attack DNNs
- Bit Flip Attack (BFA): flip a few bits among the parameters of a DNN
- Fault injection attacks (e.g., RowHammer) can induce precise bit flips
- Sparse DNNs** storing their parameters in sparse matrix formats, which can be exploited by BFAs
- Bit flip in the coordinates of nonzero weights results in “rewiring” of connections, although it does not change the value itself



Algorithm: SparseBFA

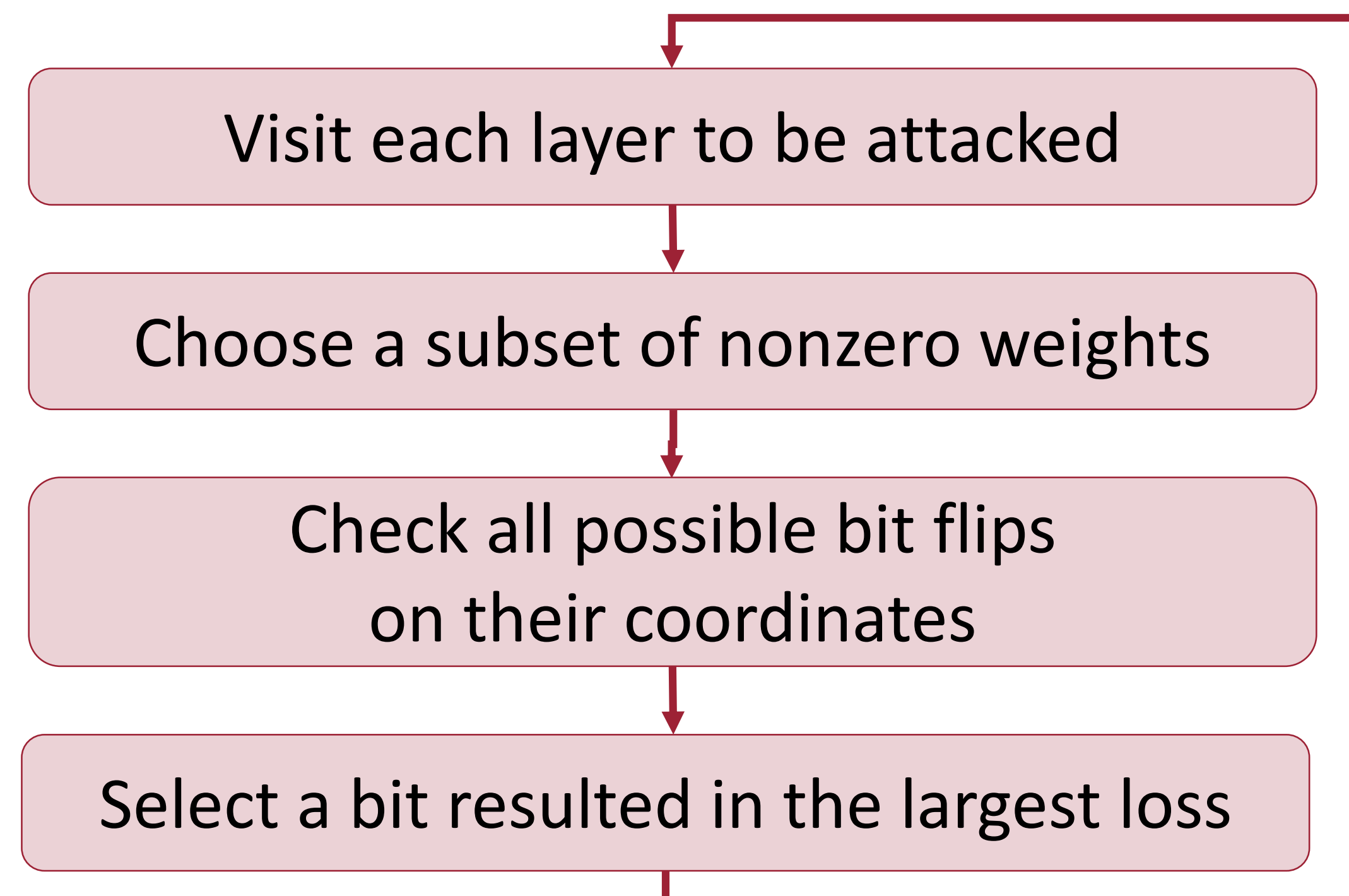
- Iteratively searches for a single bit in the coordinate lists that maximally increase the loss when it is flipped

Algorithm 1 SparseBFA for the COO format

```

1:  $m \leftarrow 0$ 
2: while  $m < M$  do  $\triangleright M$ : the maximum number of iterations
3:   loss_dict  $\leftarrow \emptyset$ 
4:   for  $i \in S_L$  do
5:      $C_i, V_i \leftarrow \text{CONVERT\_SPARSE}(W_i)$ 
6:      $S_j \leftarrow \text{GET\_CANDIDATES\_INDEX}(C_i)$ 
7:     for  $j \in S_j$  do
8:       for  $k \leftarrow 0, d$  do
9:         for  $n \leftarrow 0, n_c$  do
10:           $C_i^* \leftarrow C_i$ 
11:           $C_i^*[j][k][n] \leftarrow 1 - C_i[j][k][n]$ 
12:          if IS_VALID( $C_i^*$ ) then
13:             $W_i^* \leftarrow \text{RECONSTRUCT}(C_i^*, V_i)$ 
14:             $\Theta^* \leftarrow \Theta \setminus \{W_i\} \cup \{W_i^*\}$ 
15:            loss  $\leftarrow l(f(x; \Theta^*), t)$ 
16:            loss_dict[(i, j, k, n)]  $\leftarrow$  loss
17:   ( $i_m, j_m, k_m, n_m$ )  $\leftarrow$  GET_MAX_KEY(loss_dict)
18:    $C_{i_m}, V_{i_m} \leftarrow \text{CONVERT\_SPARSE}(W_{i_m})$ 
19:    $C_{i_m}[j_m][k_m][n_m] \leftarrow 1 - C_{i_m}[j_m][k_m][n_m]$ 
20:    $W_{i_m} \leftarrow \text{RECONSTRUCT}(C_{i_m}, V_{i_m})$ 
21:    $m \leftarrow m + 1$ 

```



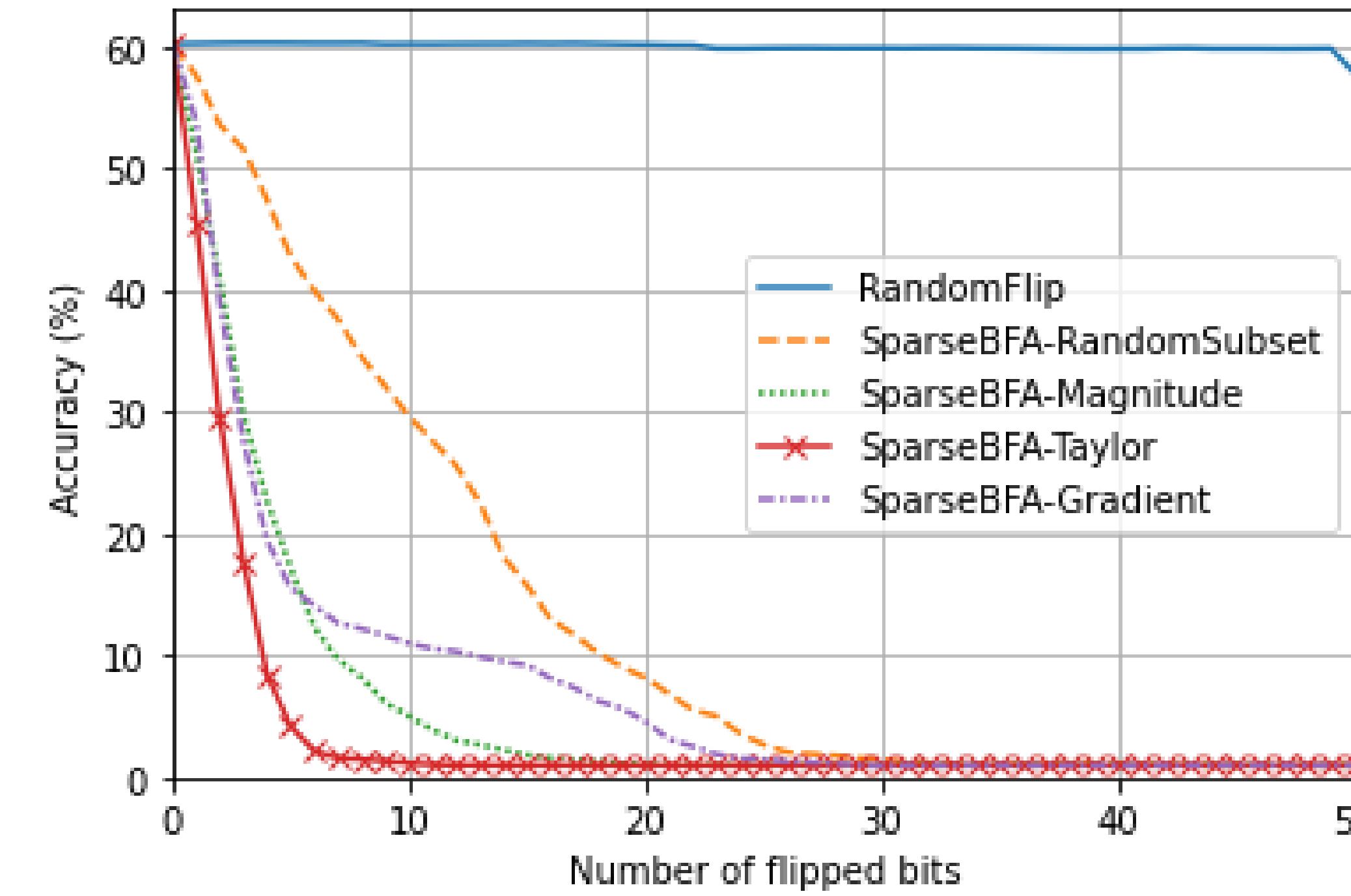
- Complexity of SparseBFA** depends on the second step for choosing a subset

$$W_i = [N_1, \dots, N_d] \xrightarrow{r: \text{sparsity}} r \cdot \prod_{i=1}^d N_i \cdot \sum_{i=1}^d \log N_i \text{ forward propagations}$$

- Choosing all nonzero weights: exhaustive search \rightarrow Exact, but slow
- Approximation methods to speed up: quantify importance of nonzero weights similar to pruning (i.e., magnitude, 1st order Taylor's)

Results

- Effectiveness of SparseBFA



- ResNet-50, TinyImageNet, excluding 1st/last layers from attack
- Random bit flips do not cause significant accuracy degradation
- SparseBFA using 1st order Taylor's term as an importance metric shows the best performance (avg. 11.8 bit flips to reach <1% top-1 accuracy)

- DNNs with various architectures and datasets are vulnerable to SparseBFA

| Dataset | Model | Sparsity | Total Bits | Original Accuracy | Target Accuracy | # Bit Flips to Reach Target | |
|--------------|----------|----------|------------|-------------------|-----------------|-----------------------------|--------------------|
| | | | | | | All Layer | Except First, Last |
| CIFAR-10 | Conv4FC2 | 98.66% | 14.26M | 86.83% | < 11% | 1.4 | 3.0 |
| | WRN | 94.37% | 8.34M | 93.14% | | 1.0 | 12.8 |
| | MobileV2 | 94.37% | 12.98M | 90.06% | | 6.8 | 8.8 |
| TinyImageNet | ResNet50 | 94.37% | 36.50M | 60.18% | < 1% (top-1) | 3.2 | 11.8 |
| | MobileV2 | 92.50% | 17.73M | 43.88% | | 1.4 | 2.8 |
| PASCAL-VOC | SSD300 | 82.20% | 137.38M | 0.70 mAP | 0.10 mAP | 34.2 | - |

- First and last layers are more critical: often a single bit flip in the first layer causes degradation to the random guess level
- Sparser DNNs are more vulnerable since SparseBFA has larger search space as sparsity grows
 - Example: Conv4FC2, CIFAR-10, excluding 1st/last layers

