

463 A Supplementary Material

464 A.1 Implementation Details

465 **Decay pruning rate with cosine annealing.** In our subspace pruning/recovery process, we let the
 466 clients prune out α_t percentage of parameters and recover the same amount of parameters to search
 467 for the subspace that fits their data. The parameter α_t is decayed with the initial rate α_0 with cosine
 468 annealing, which can be formalized as follows:

$$\alpha_t = 0.5 \times \alpha_0 \times \left(1 + \cos \left(\frac{t}{T_{end}} \pi \right) \right) \quad (7)$$

469 where t is the number of communication round, and T_{end} is the round that the mask searching is
 470 ended (notice that $\alpha_{T_{end}} = 0$). In our implementation, we set $T_{end} = T$.

471 **ERK sparsity initialization.** We use Erdős–Rényi Kernel (ERK) (Evcı et al., 2020), an empirical
 472 sparsity distribution technique, to distribute sparsity to different layers of a model. Specifically, the
 473 active parameters of the convolutional layer initialized by ERK are proportional to $1 \frac{n^{l-1} + n^l + w^l + h^l}{n^{l-1} * n^l * w^l * h^l}$,
 474 where n^{l-1} , n_l , w^l and h^l respectively specify the number of input channels, output channels and
 475 kernel’s width and height in the l -th layer. For the linear layer, the number of active parameters scale
 476 with $1 \frac{n^{l-1} + n^l}{n^{l-1} * n^l}$ where n^{l-1} and n^l are the number of neurons in the $(l-1)$ -th and l -th layer. ERK
 477 initialization, in essence, gives more sparsity to those layers with a larger number of parameters while
 478 pruning less on those small layers.

479 **Subspace pruning.** In the mask searching process, we use parameter’s magnitude to guide the
 480 pruning of model parameters. We present the PyTorch style code in Algorithm 2 to illustrate the
 481 pruning process and, correspondingly, the update of mask. Note that we only prune out parameters
 482 that are within the current subspace. Therefore, in line 6, we set the parameters that are out of
 483 the subspace to a very large value to prevent from selecting them. After that, we filter out those
 484 parameters with the smallest α_t percentage of magnitude and prune them out of the subspace.

Algorithm 2 PyTorch style code for pruning and recovery

```

1: function Prune_subspace(  $\alpha_t, \mathbf{w}_{i,t,K}, \mathbf{m}_{i,t+\frac{1}{2}}$  )
2:   Init layer sparsity  $\{s_l\}$  given overall sparsity  $s$  with ERK
3:    $\mathbf{m}_{i,t+1} = \mathbf{m}_{i,t+\frac{1}{2}}$ 
4:   for  $l = 0, 1, \dots, L - 1$  do
5:      $num_{prune} = \alpha_t \times \#$  of params in the  $l$ -th layer
6:      $sort = \text{torch.where}(\mathbf{m}_{i,t}^{(l)} == 1, \text{torch.abs}(\mathbf{w}_{i,t,K}^{(l)}), 1000 \times \text{torch.ones\_like}(\mathbf{w}_{i,t,K}^{(l)}))$ 
7:      $\_, idx = \text{torch.sort}((sort).view(-1))$ 
8:      $\mathbf{m}_{i,t+1}^{(l)}.view(-1)[idx[: num_{prune}]] = 0$ 
9:   end for
10:  Return  $\mathbf{m}_{i,t+1}$ 
11: end function
12:
13: function Recover_subspace(  $\alpha_t, \mathbf{w}_{i,t,0}, \mathbf{m}_{i,t}$  )
14:  Derive gradient  $\nabla f_i(\mathbf{w}_{i,t,0} \odot \mathbf{m}_{i,t})$  with one pass of local data
15:   $\mathbf{m}_{i,t+\frac{1}{2}} = \mathbf{m}_{i,t}$ 
16:  for  $l = 0, 1, \dots, L - 1$  do
17:     $num_{prune} = \alpha_t \times \#$  of params in the  $l$ -th layer
18:     $sort = \text{torch.where}(\mathbf{m}_{i,t+\frac{1}{2}}^{(l)} = 0, \text{torch.abs}(\nabla f_i^{(l)}(\mathbf{w}_{i,t,0} \odot \mathbf{m}_{i,t})), -1000 \times \text{torch.ones\_like}(\mathbf{w}_{i,t,K}^{(l)}))$ 
19:     $\_, idx = \text{torch.sort}((sort).view(-1), \text{descending}=\text{True})$ 
20:     $\mathbf{m}_{i,t+\frac{1}{2}}^{(l)}.view(-1)[idx[: num_{prune}]] = 1$ 
21:  end for
22:  Return  $\mathbf{m}_{i,t+\frac{1}{2}}$ 
23: end function

```

485 **Subspace recovery.** After pruning and before the next round training, we recover the same amount
 486 of parameters to explore other parameters outside the subspace. Following (Evcı et al., 2020), we use
 487 gradient information of the pruned model to guide the recovery process. Here we only recover the



Figure 7: Examples of BadNet, DBA, and Sinusoidal attack. Labels of poison samples are manipulated to the target label (e.g., a horse).

488 parameters out of the current subspace, and therefore we set the *sort_value* of the parameters within
 489 the current subspace to a sufficiently small value, as shown in Algorithm 2. Subsequently, we sort in
 490 descending to obtain the parameters with the largest- α_t percentage gradient magnitude, and recover
 491 them by updating masks.

492 A.2 Attack Methods

Table 10: Application of attack methods in threat models. "✓" corresponds to be applicable while "✗" corresponds to be not applicable.

Attack methods	Threat models		
	weak	medium	strong
BadNet	✓	✓	✓
DBA	✓	✓	✓
Sinusoidal	✓	✓	✓
Scaling	✗	✓	✓
FixMask (adaptive)	✗	✓	✓
Neurotoxin	✗	✗	✓
Omniscience (adaptive)	✗	✗	✓

493 As we mention in the main body, we classify the attack model into data-level attack and algorithm-
 494 level backdoor. We in the following give brief description of each data-level backdoor that we
 495 simulate in federated learning setting.

- 496 • **BadNet.** BadNet is the earliest, and also the simplest backdoor attack first proposed in (Gu et al.,
 497 2017). To perform BadNet attack, the malicious client simply add the same backdoor trigger on
 498 some of the data samples, and modify the label of these poisoned samples to the target label. In test
 499 time, the malicious clients can place the backdoor trigger on the test samples, such that the victim
 500 model can produce the target output no matter what the original test samples are.
- 501 • **DBA.** DBA (Xie et al., 2019) is a backdoor attack specifically targeted on FL. To perform DBA
 502 attack, the authors decompose the backdoor trigger into several local pattern, and assign the local
 503 pattern to different clients to poison their local data. For test time, the attacker will interpose the
 504 completed trigger on top of the test samples they want to manipulate. It is suggested by the authors
 505 that DBA is substantially more persistent and stealthy against FL. In our simulation, we decompose
 506 the "plus" trigger into 4 local patterns, and let each malicious client to be assigned each local
 507 pattern.
- 508 • **Sinusoidal attack.** Sinusoidal attack (Barni et al., 2019) shares a similar perspective with BadNet,
 509 which also utilize the same trigger for all the malicious clients to poison their samples. However,
 510 the backdoor trigger they use is a horizontal sinusoidal signal defined by $v(i, j) = \Delta \sin(2\pi j f / m)$,
 511 $1 \leq j \leq m, 1 \leq i \leq l$, for a certain frequency f . The authors claim that this design of trigger i)
 512 is weak enough to ensure the stealthiness of the attack, but also ii) be detectable in the same (or
 513 similar) feature space used by the network to classify the pristine samples. In our simulation, we
 514 adopt the default hyper-parameter $\delta = 20$ and $f = 6$ for performing this attack.

515 Examples of these data-level attacks are visually shown in Figure 7.

516 In the following we give brief description on the algorithm-level backdoor that has been simulated in
517 this paper.

- 518 • **Scaling**. The basic idea of Scaling (Bagdasaryan et al., 2020) is to enlarge the gradient update
519 when a malicious client return its update to server. This mechanism allows the malicious client to
520 enlarge its gradient’s impact on the global model, and therefore is effective when the poison ratio
521 and attacker number are small.
- 522 • **FixMask**. FixMask is an adaptive attack method specifically targeting Lockdown. In Lockdown,
523 the malicious clients are assumed to faithfully search for their subspace using their local data. For
524 FixMask attack, the malicious clients freeze their mask to be the initial mask that is shared by all
525 the clients in round $t = 0$, and refuse to change afterwards.

526 Particularly, we want to emphasize that the data-level and algorithm-level backdoor can potentially
527 be combined together to produce better attack performance. However, since this paper focus on the
528 defense aspect, we leave a more thorough study of the attack model future work. We also include two
529 advanced attack algorithms that can only be conducted given extra server information in addition to
530 permission on manipulation of the attacker’s own training process and data.

- 531 • **Neurotoxin**. Neurotoxin proposed in (Zhang et al., 2022) explores a durable attack method in
532 the scenario that the attackers can only participate limited rounds. Their main observations in the
533 limited participation case are that i) the benign update can recover the global model after attacker
534 ceases attack. ii) the majority of the l2 norm of the aggregated benign update is contained in a small
535 number of coordinates (Let’s call these benign coordinates). Utilizing the above observations, the
536 authors propose Neurotoxin, which is to let the malicious clients project their gradient update to the
537 subspace excluding the global coordinates. By this means, the projected updates from the malicious
538 clients are mostly embedded to the coordinates that have less perturbation by the benign updates
539 (which focus on the benign coordinates) after ceasing attack. However, Neurotoxin cannot escape
540 Lockdown defense in principle. There are mainly two reasons. i) Lockdown only broadcast to the
541 clients some coordinates weights (equivalently, some coordinates of gradient update) as per their
542 subspace. Therefore, Neurotoxin cannot obtain the top-k coordinate of the server gradient as benign
543 coordinates. ii) Lockdown requires clients to report the subspace that they want to update, and the
544 subspace that are substantially different from others will be pruned afterwards. In other words, if
545 the attackers adopt neurotoxin to choose the subspace that excludes the benign coordinates, their
546 subspace can be easily identified by comparing with other benign client’s subspace, and therefore
547 will be pruned out. In our simulation, we assume Neurotoxin can acquire the server gradient update
548 by some means. Therefore, it is classified as an attack method for *strong* threat model. In our
549 simulation, we set its hyper-parameter mask ratio to be 0.25.

- 550 • **Omniscience**. This is an adaptive attack that assumes the knowledge of Lockdown and try to break
551 it. The main idea is to assume the client’s has knowledge of the consensus subspace after going
552 through consensus fusion, and project their gradient update into this subspace. This efficiently
553 avoids the malicious weights to be pruned out by the consensus fusion operation. However,
554 the requirement of conducting this attack is very stringent. The malicious client needs to have
555 knowledge of the consensus subspace, which either is leaked from server, or is computed if other
556 clients’ subspace is known by the attacker. Neither of this condition is easy to establish for an
557 attacker in a federated learning system.

558 In summary, we show in Table 10 the attack methods we can perform with specific threat models.

559 A.3 Defense Methods

560 In this section, we give a brief description of the defense baseline we compare against.

- 561 • **RLR**. RLR proposed in (Ozdayi et al., 2021) utilizes coordinate-wise server learning rate to inverse
562 the gradient coordinates in which different clients have different sign. Their observation is that the
563 malicious coordinates tends to be those coordinates that have conflicting sign in gradient while
564 for the benign coordinates that are not poisoned, most of the clients will agree with their sign.
565 Therefore, by looking at the gradient update from clients, the server is able to identify the malicious
566 coordinates and subsequently inverse its sign in the aggregation phase. However, the malicious
567 clients are able to launch adaptive attack if he knows the gradient update downloaded from server.

- 568 • **RFA.** Aiming at defense against corrupted updates from clients, RFA (Pillutla et al., 2022) utilizes
569 the concept of geometric medium to aggregate the gradient update from clients. Geometric medium
570 avoids the gradient that has excessively large norm (usually is the malicious one) to impact too
571 much on the averaging process. Specifically, when doing aggregation, instead of directly averaging
572 the uploaded gradient, the server aims to obtain global model v that minimizes: $\sum_{i=1}^m \|v - w_i\|$,
573 and w_i is the uploaded local model. This problem is solved by the Smoothed Weiszfeld Algorithm.
574 Similar techniques are studied in (Sifaou & Li, 2022), (Ghosh et al., 2019) and (Cao et al., 2020).
- 575 • **Krum.** Targeting Byzantine attack, Krum (Blanchard et al., 2017) adopts the idea of finding
576 the gradient update that is closest to its $n - f - 2$ neighbours such that it can ensure (α, f) -
577 Byzantine resilience where α is the angle depends on the ratio of the deviation over the gradient
578 , f is the number of attackers. Specifically, Krum aims to find the the i^* -th client that minimize
579 $s(i) = \sum_{i \rightarrow j} \|V_i - V_j\|^2$ where $i \rightarrow j$ denotes the set of i 'th client's $n - f - 2$ closest neighbours,
580 and V_i denotes the gradient update from client i . After identifying i^* , Krum returns V_{i^*} as the
581 robust gradient used for aggregation.
- 582 • **Trimmed mean.** Trimmend mean is proposed in (Yin et al., 2018) to counter byzantine failures in
583 the distributed machine learning scenario. Their high level idea is to exclude the outlier gradient
584 value when doing aggregation. Specifically, before aggregation, the server coordinate-wise trims
585 the TopK gradient and the bottomK gradient among those uploaded gradient. After trimming,
586 the server assume the outlier has been trimmed, and directly average the clean gradient. In our
587 simulation, we set the trim ratio to be 0.1.

588 A.4 Security Analysis

589 We make the following observations on Lockdown's security performance. Lockdown can success-
590 fully defend all the data-level attack, i.e., the attack falls in to the scope of *weak* threat models. For
591 the algorithm-level attacks, we have incorporated an adaptive attack targeting on Lockdown and a
592 gradient scaling method into study. Our results show that Lockdown can also defend all the attacks we
593 have tested. However, since the algorithm-level attacks are more adaptive, we cannot make guarantee
594 that Lockdown is unbreakable by any algorithm-level attacks, especially those that are specially
595 designed for Lockdown. For advanced attack that allows attacker to acquire server's information, we
596 create another adaptive attack Omniscience that can successfully circumvent Lockdown's defense.
597 Performing Omniscience attack needs the attacker to know about the consensus subspace. However, it
598 is challenging, if not impossible, for the attacker to infer the consensus subspace, since only a subset
599 of the server gradient update is distributed to clients, further constraining the global information
600 access of the attackers.

601 A.5 More Visualization

602 **Input-level visualization.** In Figure 8, we add additional experiments to visualize the gradient w.r.t
603 the input of the first layer, which visually explains how different semantic information within the
input image contributes to activating the target output neuron.

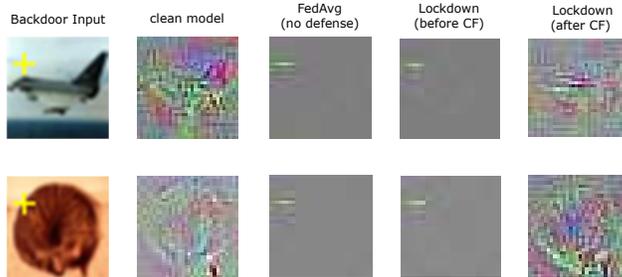


Figure 8: Smooth grad (Smilkov et al., 2017) visualization of models given backdoor input. The first column is data input with backdoor trigger. The subsequent columns demonstrate the gradient with respect to the input of i) a model without being poisoned. ii) a model trained by FedAvg with poisoned data, iii) Lockdown's global model under poisoning before going through consensus fusion (CF) and iv) Lockdown's final model. A clean model emphasize the correct semantic within the input, e.g., wing of a plane, while a poisoned model emphasizes the yellow "plus" backdoor trigger.

604 **Parameters-level visualization.** In Figure 9, we visualize the projected parameters produced by
 605 Lockdown. The experiment is conducted on MNIST with a two-layer MLP model. After reducing its
 606 output dimension and reshaping it into the original input, we plot the projected absolute weights of the
 607 first layer of MLP. As found, by projecting the global weights into malicious client’s subspace (left),
 608 the corresponding connectivity that joint the backdoor trigger still present. However, by projecting
 609 the global weights into one of the benign client’s subspace (middle), the backdoor trigger no longer
 610 connects with large absolute weights. The same phenomenon is observed for the consensus subspace
 611 after going through consensus fusion (right).

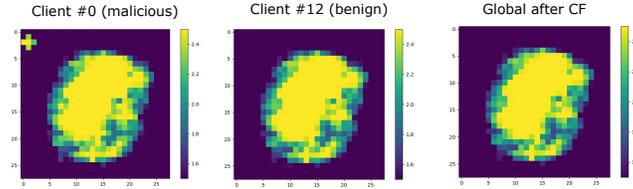


Figure 9: Visualization of absolute global weights after projecting into the local or global subspace. Left: projecting into local subspace of a malicious client. Middle: projecting into local subspace of a benign client. Right: projecting into consensus subspace produced by consensus fusion. The brighter the color is, the feature locates in that part is more important. The bright backdoor trigger "+" is not visible in the middle and right image. See more details in the main text.

612 A.6 Ablation study

613 We perform ablation study of Lockdown on CIFAR10. BadNet is the default attack method.

614 **Gradient-based recovery vs. random recovery.** In subspace recovery process, we use gradient
 615 magnitude to guide the recovery of parameters. In Table 11, we show the empirical comparison
 616 between the gradient-based recovery and random recovery. The results showcase that recovery with
 617 the gradient can significantly reduce the ASR (by up-to 78.3% reduction) though the benign acc of
 618 the model suffered a little bit (by up-to 2.3% drop). This is because gradient magnitude tends to
 619 guide the subspace searching process to acquire heterogeneous subspaces for clients with different
 620 training data. With more heterogeneous subspaces, the knowledge transferring between clients will
 621 be deterred since their the subspace overlap is small, which leads to the degradation of benign acc. On
 622 the other hand, small subspace overlap can also facilitate the process of de-poisoning by consensus
 623 fusion, which leads to a reduction of ASR.

Table 11: Ablation study for parameters recovery implementation.

Methods (IID)	Benign Acc(%) ↑	ASR(%) ↓	Backdoor Acc(%) ↑
Random recovery	91.1	13.0	79.7
Recovery w/ gradient (ours)	90.7	1.4	87.8
Methods (Non-IID)	Benign Acc(%) ↑	ASR(%) ↓	Backdoor Acc(%) ↑
Random recovery	88.9	17.2	74.3
Recovery w/ gradient (ours)	84.9	2.2	78.4

624 **ERK initialization vs. uniform initialization.** In the SubspaceInit() function, we use ERK to
 625 allocate the sparsity of each layer in a model. To justify the necessity of ERK initialization, we
 626 replace the ERK initialization with uniform initialization, which uniformly allocates sparsity to each
 627 layer. As shown in Table 12, uniform initialization will largely compromise the benign accuracy
 628 and slightly increase the ASR. This justifies that the sparsity should be set larger for the layer with a
 629 larger number of parameters (which essentially is what ERK does).

630 **Consensus fusion (CF).** In Figure 10, we demonstrate the necessity of consensus fusion under
 631 different poison ratios. With consensus fusion, benign accuracy is significantly increased by up-to 60%
 632 while the ASR is reduced by up-to 80%. This result shows that masking out some malicious/dummy
 633 parameters can perturb the backdoor function and thereby curing the poisoned model.

Table 12: Ablation study for sparsity initialization.

Methods (IID)	Benign Acc \uparrow	ASR \downarrow	Backdoor Acc \uparrow
Uniform	81.7	5.7	75.9
ERK (ours)	90.1	7.1	83.7
Methods (Non-IID)	Benign Acc \uparrow	ASR \downarrow	Backdoor Acc \uparrow
Uniform	75.5	3.1	70.6
ERK (ours)	86.1	3.4	82.2

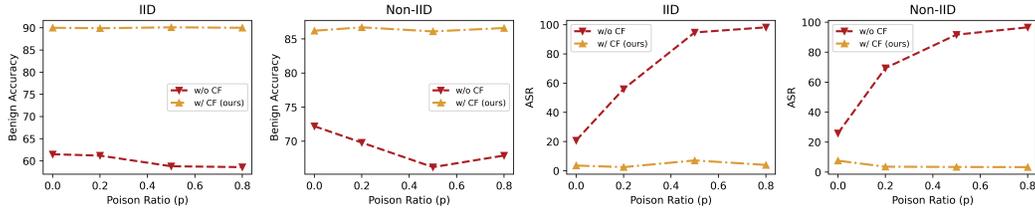


Figure 10: Impact of consensus fusion of Lockdown.

634 **A.7 Hyper-parameter Sensitivity Analysis**

635 In this section, we perform hyper-parameter sensitivity analysis for lockdown. The evaluation is
 636 conducted on CIFAR10 under the default simulation setting in Table 2 unless otherwise specified.

637 **Sparsity s .** In Table 13, we set other hyper-parameters as default and tune the sparsity to different
 638 levels. As shown, Lockdown loses its defense efficacy when sparsity is low. This phenomenon is
 639 understandable since Lockdown reduces to FedAvg when sparsity is 0. On the other hand, with
 640 larger sparsity, the benign accuracy of the model suffers due to the reduction of trainable parameters.
 641 Therefore, there exists a tradeoff for the sparsity of Lockdown. Larger sparsity promises lower model
 complexity, smaller comm overhead, and also lower ASR, but at the cost of losing benign accuracy.

Table 13: Performance of Lockdown under different sparsity s .

s (IID)	Benign Acc \uparrow	ASR \downarrow	# of params \downarrow
0	91.0	68.4	6.57M
0.2	90.9	61.1	5.26M
0.5	91.0	10.9	3.29M
0.75	90.1	7.1	1.65M
0.9	88.3	3.0	0.66M
s (Non-IID)	Benign Acc \uparrow	ASR \downarrow	# of params \downarrow
0	89.1	70.3	6.57M
0.2	88.4	52.6	5.26M
0.5	87.1	14.1	3.29M
0.75	86.1	3.4	1.65M
0.9	85.0	2.9	0.66M

642

643 **Initial pruning/recovery rate.** We also show the effect of initial pruning/recovery for the learning
 644 performance. As shown, larger pruning rate would typically results in the drop of benign accuracy
 645 but also enhance the ASR under poisoning attack. Specially, when $a_0 = 0$, lockdown reduces to
 646 train a sparse subnetwork from scratch, without evolving the sparse coordinate. This setting cannot
 647 eliminate the "poison-couple" effect, therefore the ASR is as high as FedAvg with no defense. On the
 648 other hand, setting a_0 will also result in isolation of subspace for different clients, resulting in lack of
 649 consensus in the global space and therefore leading to drop of benign accuracy.

650 **Consensus fusion threshold θ .** In Figure 11, we tune the CF threshold θ to see its impact on
 651 different settings of attacker number N . In all settings of N , we see that: i) θ should not be set to
 652 be too small; otherwise, the benign accuracy would be lower, and the ASR will be higher. ii) θ also
 653 should not be set too large; otherwise, it will severely compromise benign accuracy, but the reduction

Table 14: Performance of Lockdown under different initial pruning/recovery rate a_0 .

a_0 (IID)	Benign Acc \uparrow	ASR \downarrow	Backdoor Acc \uparrow
0	90.5	49.4	47.7
1e-5	90.5	5.2	84.3
1e-4	90.1	7.1	83.7
1e-3	88.1	3.7	84.7
1e-2	87.2	3.5	83.7
1e-1	87.0	3.1	83.4

a_0 (Non-IID)	Benign Acc \uparrow	ASR \downarrow	Backdoor Acc \uparrow
0	88.5	85.3	14.0
1e-5	87.4	8.5	78.8
1e-4	86.1	3.4	82.2
1e-3	84.9	2.1	80.4
1e-2	83.4	5.3	76.4
1e-1	83.7	5.2	77.6

654 of ASR will not be too significant. Per our results, the consensus threshold should be chosen carefully
 655 according to the number of attackers, which of course, is unknown in most cases. However, given that
 656 the attackers within the system should not take up a large portion, θ set to be 50% of the total number
 of clients will be sufficient to counteract the effect of backdoor attack in a general attack scenario.

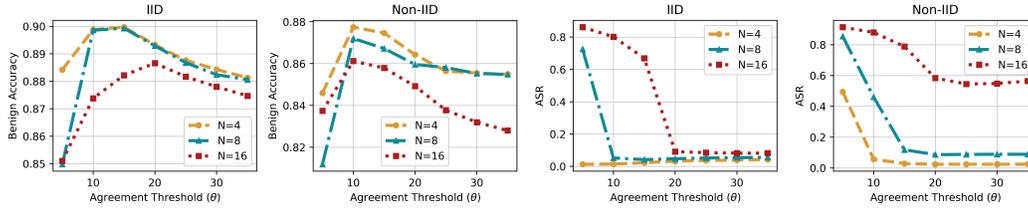


Figure 11: Impact of consensus fusion threshold of Lockdown in different # of attackers setting.

657

658 A.8 Limitations

659 Our method utilizes sparsity of model to counter backdoor attack. However, we are aware that sparsity
 660 in its current stage can hardly guarantee acceleration of the training/inference speed. At present, the
 661 current sparse acceleration technique requires 2:4 sparse operation. More specifically, the 2:4 sparse
 662 operation requires that there are at most two non-zero values in four contiguous memory, which
 663 may not hold for the sparse model produced by Lockdown. But we insist that our method has great
 664 potential to achieve truly training acceleration with development of sparse technique.

665 There are potentially other adaptive backdoor attacks that can break the defense of lockdown,
 666 especially under the assumption that attackers have full control over its local training process and
 667 has knowledge of the defense. We leave the research of potential attacks against Lockdown as future
 668 works.

669 A.9 Broader Impact

670 The poison-coupling effect we discover in this paper might be mis-used to guide the design of
 671 backdoor attack method in centralized learning/FL scenario. We will continue this line of research
 672 and further propose attack/defense method to better study/mitigate such an effect. We also open-
 673 source our code to facilitate researchers/machine learning engineer in academy/industry to study and
 674 understand the discovered phenomenon.