

---

# Supplementary Material of “BackdoorBench: A Comprehensive Benchmark of Backdoor Learning”

---

Baoyuan Wu<sup>1\*</sup> Hongrui Chen<sup>1</sup> Mingda Zhang<sup>1</sup> Zihao Zhu<sup>1</sup>  
Shaokui Wei<sup>1</sup> Danni Yuan<sup>1</sup> Chao Shen<sup>2</sup>

<sup>1</sup>School of Data Science, Shenzhen Research Institute of Big Data,  
The Chinese University of Hong Kong, Shenzhen

<sup>2</sup>School of Cyber Science and Engineering, Xi’an Jiaotong University

## A Additional information of backdoor attack and defense algorithms

### A.1 Descriptions of backdoor attack algorithms

In addition to the basic information in Table 1 of the main manuscript, here we describe the general idea of eight implemented backdoor attack algorithms in BackdoorBench, as follows.

- **BadNets** [6]: It was the first work in backdoor learning, which simply inserted a small patch with fixed pattern and location to replace the original pixels in the clean image to obtain a poisoned image.
- **Blended** [3]: It extended BadNets by encouraging the invisibility of the trigger through alpha blending.
- **Label consistent (LC)** [20]: It generated a poisoned image using adversarial attack, by enforcing it to be close to the clean target image in the original RGB space, and close to the clean source image patched with a trigger in the feature space of a pre-trained clean model. Since the poisoned image is labeled as the target class, the mapping from the trigger to the target class could be learned.
- **SIG** [1]: It adopted a sinusoidal signal as the trigger to perturb the clean images of the target class, while not changing their labels, such that achieving the label consistent backdoor attack.
- **Low frequency attack (LF)** [30]: It was built upon an analysis that the triggers in many backdoor attacks bring in high-frequency artifacts, which are easily detectable. Inspired by this analysis, LF developed a smooth trigger by filtering high-frequency artifacts from a universal adversarial perturbation.
- **Sample-specific backdoor attack (SSBA)** [10]: It utilized an auto-encoder to fuse a trigger (*e.g.*, a string) into clean samples to obtain poisoned samples. The residual between the poisoned and the clean sample varied for different clean images, *i.e.*, sample-specific.
- **Input-aware dynamic backdoor attack (Input-aware)** [14]: It was a training-controllable attack by simultaneously learning the model parameters and a trigger generator. When testing, the learned trigger generator generated one unique trigger for each clean testing sample.
- **Warping-based poisoned networks (WaNet)** [15]: It was also a training-controllable attack. A fixed warping function is adopted to slightly distort the clean sample to construct the poisoned sample. The attacker further controlled the training process to ensure that only the adopted fixed warping function can activate the backdoor.

---

\*Corresponds to Baoyuan Wu (wubaoyuan@cuhk.edu.cn).

## A.2 Descriptions of backdoor defense algorithms

In addition to the basic information in Table 2 of the main manuscript, here we describe the general idea of nine implemented backdoor defense algorithms in BackdoorBench, as follows.

- **Fine-tuning (FT)**: It is assumed that fine-tuning the backdoored model on a subset of clean samples could mitigate the backdoor effect. Note that FT is a widely used approach for transferring pre-trained models to new tasks, but it has been used as a basic component in several backdoor defense methods, such as Fine-pruning (FP) [12], Neural Attention Distillation (NAD) [8].
- **Fine-pruning (FP)** [12]: It is built upon the assumption that *poisoned and benign samples have different activation paths in the backdoored model*. Inspired, FP proposed to firstly prune some inactivated neurons of clean samples, then fine-tune the pruned model based on the subset of benign samples to recover the model performance.
- **Neural Attention Distillation (NAD)** [8]: Its assumption is same with FT. Instead of directly use the fine-tuned model as the mitigated model, NAD adopts the first fine-tuned model as a teacher, and fine-tunes the backdoored model again by encouraging the consistency of the attention representation between the new fine-tuned model and the teacher model.
- **Neural cleanse (NC)** [26]: It is built upon the assumption that *the trigger provides a “shortcut” between the samples from different source classes and the target class*. Based on this assumption, the possible trigger is searched through optimization. If a small-size trigger (e.g., a small patch in the image) is found, then the model is detected as backdoored model, which is then mitigated through pruning based on the searched trigger.
- **Adversarial Neuron Pruning (ANP)** [27]: It is built upon an observation that *the neurons related to the injected backdoor are more sensitive to adversarial neuron perturbation (i.e., perturbing the neuron weight to achieve adversarial attack) than other neurons in a backdoored model*. Inspired by this, ANP proposed to prune these sensitive neurons for backdoor mitigation.
- **Activation Clustering (AC)** [2]: It is built upon an observation that *the sample activations (i.e., the feature presentations) of the target class will form two clusters, and the smaller cluster corresponds to poisoned samples, while those of other classes form one cluster*. Then, the model is trained from scratch based on the dataset without poisoned samples.
- **Spectral Signatures (SS)** [24]: Its assumption is that *the feature representation distributions of benign and poisoned samples in one class are spectrally separable*, which is a concept of robust statistics. Consequently, the poisoned samples can be identified through analyzing the spectrum of the covariance matrix of the feature representations. Then, the model is retrained from scratch by removing the poisoned samples from the training set.
- **Anti-Backdoor Learning (ABL)** [9]: It is built upon an observation that *the loss values of poisoned samples drops much faster than those of benign samples in early epochs during the training process*. Inspired, ABL proposed to firstly isolate poisoned samples from benign samples according to their difference on loss dropping speed, then mitigate the backdoor effect by maximizing the loss of the isolated poisoned samples.
- **Decoupling-based Backdoor Defense (DBD)** [7]: It is built upon an observation that *poisoned samples from different samples will gather together in the feature space of a backdoored model*. DBD proposed to prevent the gathering by learning the model backbone through self-supervised learning without labels, rather than the standard supervised learning. Then, since the poisoned samples are separated, their loss values are larger than benign samples when learning the classifier, such that samples with large loss values can be identified as poisoned samples. Finally, the labels of poisoned samples are abandoned, and a semi-supervised fine-tuning of both the backbone and classifier is conducted to improve the model performance.

---

<sup>1</sup>For FP and ANP, we define a hyper-parameter *the tolerance of clean accuracy reduction* as the maximum relative drop of clean accuracy. It is used to determine the number of pruned neurons.

<sup>2</sup>For AC and Spectral,  $N_{fe}$  is the dimensions of the representation.

<sup>3</sup>For NAD, we use the code and recommended hyper-parameters at <https://github.com/bboylyg/NAD/tree/d61e4d74ee697f125336bfc42a03c707679071a6>.

Table 1: Hyper-parameter settings of all implemented attack methods.

Attack	Hyper-parameter	Setting	Theoretical complexity
General Settings	attack target	all-to-one with class 0	train sample size ( $N_{\mathcal{T}}$ ),
	optimizer	SGD	batch size ( $B$ ),
	momentum	0.9	forward process for net C ( $f_C$ ),
	weight decay	0.0005	backward process for net C ( $b_C$ ),
	batch size	128	epochs ( $E$ ),
	epochs on CIFAR10 and CIFAR-100	100	number of classes ( $N_{cls}$ )
	epochs on GTSRB	50	
	lr schedule (except for training-controllable attack) on CIFAR10, CIFAR-100 and GTSRB	CosineAnnealingLR	
epochs on TinyImageNet-200	200		
lr schedule on TinyImageNet-200	ReduceLRonPlateau		
random seed	0		
BadNets [6]	pattern & location	$3 \times 3$ , pure white, at downright corner (no margin left)	$O(N_{\mathcal{T}}/BE(f+b))$
Blended [3]	pattern	hello kitty	$O(N_{\mathcal{T}}/BE(f+b))$
	alpha	0.2	
Label Consistent [25]	adversarial attack	PGD	$O(N_{\mathcal{T}}/BE(f+b) + N_{\mathcal{T}}/BT_{step}(f+b))$
	step ( $T_{step}$ )	100	
	$\alpha$	1.5	
SIG [1]	$\epsilon$	8	$O(N_{\mathcal{T}}/BE(f+b))$
	$\delta$	40	
Low Frequency [30]	$f$	6	$O(T_{ter}N_{sample}(2f + T_{df}f + T_{df}N_{cls}b) + N_{\mathcal{T}}/BE(f+b))$
	maximum number of termination iteration ( $T_{ter}$ )	50	
	fooling rate	0.2	
	overshoot	0.02	
	maximum number of iterations for deepfool ( $T_{df}$ )	200	
sample number for UAP ( $N_{sample}$ )	100		
SSBA [10]	encoded bit	1	$O(T_{step}(f_{auto} + b_{auto}) + N_{\mathcal{T}}/BE(f+b))$
	autoencoder train step ( $T_{step}$ )	140000	
Input-aware [14]	Generator lr (both for M, G)	0.01	$O(E_{mask}N_{\mathcal{T}}/B(f_{mask} + b_{mask}) + (E - E_{mask})N_{\mathcal{T}}/B(2f_{mask} + 2f_{generator} + 2b_{generator} + f + b))$
	schedule (for M, G and C)	MultiStepLR	
	schedule milestones for G	200, 300, 400, 500	
	schedule milestones for C	100, 200, 300, 400	
	schedule milestones for M	10, 20	
	schedule gamma (for M and G)	0.1	
	$\lambda_{div}$	1	
	$\lambda_{norm}$	400	
	mask density	0.032	
	cross_ratio	1	
mask train epochs ( $E_{mask}$ )	25		
WaNet [15]	cross_ratio	2	$O(N_{\mathcal{T}}/BE(f+b))$
	lr schedule	MultiStepLR	
	schedule milestones	100, 200, 300, 400	
	grid_rescale	1	

### A.3 Implementation details and computational complexities

**Running environments** Our evaluations are conducted on GPU servers with 2 Intel(R) Xeon(R) Platinum 8170 CPU @ 2.10GHz, RTX3090 GPU (32GB) and 320 GB RAM (2666MHz). The versions of all involved softwares/packages are clearly described in the README file of the Github repository (see <https://github.com/SCLBD/BackdoorBench>). Here we didn't repeat the descriptions.

**Hyper-parameter settings** The hyper-parameter settings adopted in our evaluations about backdoor attack and defense algorithms are described in Table 1 and Table 2, respectively. With these hyper-

Table 2: Hyper-parameter settings of all implemented defense methods.

Defense	Hyper-parameter	Setting	Theoretical complexity
General Settings	optimizer	SGD	train sample size ( $N_T$ ), batch size ( $B$ ), forward process for net C ( $f_C$ ), backward process for net C ( $b_C$ ), epochs ( $E$ ), number of classes ( $N_{cls}$ ) the number of pruning neurons ( $N_{neu}$ )
	momentum	0.9	
	weight decay	0.0005	
	batch size	256	
	epochs on CIFAR10, CIFAR-100 and GTSRB	100	
	lr schedule (except for special learning defense) on CIFAR10, CIFAR-100 and GTSRB	CosineAnnealingLR	
	epochs on TinyImageNet-200	200	
lr schedule on TinyImageNet-200	ReduceLRonPlateau		
the number of pruning neurons	the number of neurons in the last layer		
random seed	0		
FT	ratio of validation data ( $p_v$ )	5%	$O(N_T p_v / BE(f + b))$
FP [12]	ratio of validation data ( $p_v$ )	5%	$O(N_{neu} N_T p_v f / B)$
	the tolerance of accuracy reduction <sup>1</sup>	10%	$+O(N_T p_v / BE(f + b))$
NAD <sup>3</sup> [8]	ratio of validation data for teacher model ( $p_v$ )	5%	$O(N_T p_v / BE_{ft}(f + b))$
	$\beta_1$ for the loss	500	$+O(N_T p_v / BE(2f + b))$
	$\beta_2$ for the loss	1000	
	$\beta_3$ for the loss	1000	
	the power for attention	2.0	
	the epoch for teacher model to fine-tune ( $E_{ft}$ )	10	
NC [26]	the norm used for the reversed trigger	L1	$O(N_T p_v E_r N_{cls} / B(f + b))$
	cleaning ratio ( $p_v$ )	0.05	$+O(N_T p_v / BE(f + b))$
	unlearning ratio	0.2	
	the epoch of learning trigger ( $E_r$ )	80	
ANP [27]	the tolerance of accuracy reduction <sup>1</sup>	10%	$O(e_i(f + b))$
	number of validation data ( $p_v$ )	5%	
	the number of iteration during perturbation ( $e_i$ )	2000	
	$\epsilon$	0.4	
	$\alpha$	0.2	
AC [2]	number of reduced dimensions	10	$O(N_T / BN_{fc}^3) + O(N_T / BE(f + b))^2$
Spectral [24]	The percentile of backdoor data	85%	$O(N_{fc}^3) + O(N_T / BE_d(f + b))^2$
ABL [9]	tuning epochs ( $E_{tu}$ ) for CIFAR10, CIFAR-100 and GTSRB	20	$O(N_T / BE_{tu}(f + b))$
	finetuning epochs ( $E_{ft}$ ) for CIFAR10, CIFAR-100 and GTSRB	60	$+O(N_T(1 - p_i) / BE_{ft}(f + b))$
	unlearning epochs ( $E_u$ ) for CIFAR10, CIFAR-100 and GTSRB	20	$+O(N_T p_i / BE_u(f + b))$
	tuning epochs for Tiny	40	
	finetuning epochs for Tiny	120	
	unlearning epochs for Tiny	4	
	lr for unlearning	0.0005	
	the value of flooding	0.5	
	the isolation ratio of training data ( $p_i$ )	0.01	
DBD [7]	the epoch for self-supervised learning ( $E_{se}$ )	100	$O(N_T / B_{se} E_{se}(f + b))$
	the epoch for warmup ( $E_{wa}$ )	10	$+O(N_T / B_{semi} E_{se}(f + b))$
	the epsilon for the dataset ( $\epsilon$ )	0.5	$+O(N_T \epsilon / B_{semi} E(f + b))$
	during the semi-supervised learning		
	The batch size of self learning ( $B_{self}$ )	512	
	The batch size of self learning ( $B_{semi}$ )	128	

parameter settings, the reported results of 8,000 pairs of evaluations could be reproduced. Moreover, we would like to explain our rules to adopt above settings, as follows:

- **We don't perform a separate hyper-parameter search for each method**, mainly due to the following two reasons:
  - As shown in Tables 1 and 2, most methods have several hyper-parameters. For most hyper-parameters of a method, there is neither a good rule to determine the values, nor a suitable range of the values suggested in its original manuscript. And, the suitable value or range of each hyper-parameter may vary across different datasets, different model architectures, different against attack/defense methods. Consequently, the hyper-parameter search space for each method could be very large, requiring unimaginably high computational resource.
  - Even assuming sufficient computing resources, then we can search a good value for each hyper-parameter of each method in each evaluation. However, the comparison results and analysis based on sufficient hyper-parameter search may be unfair and make no sense in practice. Because, we still cannot tell a rule or even some experiences to determine the hyper-parameter values in practice. The sensitivity to hyper-parameters should also be an important metric of one method's performance, not just the best ACC/ASR values through the sufficient hyper-parameter search.
- **How do we set the hyper-parameter values in our current 8000 pairs of evaluations.**

- If the original paper has provided the suggested good values of some hyper-parameters, then we adopt those values in our evaluations. For example, the ANP defense method explicitly wrote that "the perturbation budget  $\epsilon = 0.4$  and the trade-off coefficient  $\alpha = 0.2$ ", so we also adopt these values in our evaluations.
- For those hyper-parameters without suggested values/ranges (or even without descriptions) in their original papers, we will search values that lead to comparable results (ACC/ASR) with the reported results in the same setting (*i.e.*, same dataset, same/similar model architecture, same poisoning ratio), then fix these values in evaluations of other settings (*e.g.*, changing the poisoning ratio).
- The consistent values of hyper-parameters of each method across different settings somewhat guarantee the fairness of evaluations. And, since the adopted values may not be the optimal ones for some hyper-parameters, we didn't conduct the fine-grained analysis about the effects of some specific hyper-parameters (*e.g.*, the trigger size/location in attack methods with patch based triggers). Instead, we provided some high-level analysis *w.r.t.* the shared hyper-parameters in all methods (*e.g.*, the number of classes, the poisoning ratio, the model architecture). The findings of these high-level analysis will not be significantly affected by the particular hyper-parameters of each individual method.

**Computational complexities** The computational complexity of each attack and each defense algorithm is also described in Table 1 and Table 2, respectively.

## B Additional evaluations and analysis

### B.1 Full results on CIFAR-10

The full results on CIFAR-10 with five different poisoning ratios (*i.e.*, 10%, 5%, 1%, 0.5%, 0.1%) are presented in Tables 10 – 14, respectively. The remaining results on other datasets and model architectures among 8,000 attack-defense pairs of evaluations are presented in the Leaderboard in the BackdoorBench website (see <https://backdoorbench.github.io>).

### B.2 Results overview

In Figure 1, we present the performance distribution of attack-defense pairs on PreAct-ResNet18 and VGG-19 with two poisoning ratios of 5% and 10%, respectively. As we mentioned in Section 4.2, if we measured the effectiveness of methods by clean accuracy (C-Acc) and ASR, a perfect attack method should be located at the top-right corner; the perfect defense method should show in the top-left corner. If we measured robust accuracy (R-Acc) and ASR, the reduced ASR value would be desirable to equal the increased R-Acc. This defense method can recover the correct prediction and eliminate the backdoor successfully. Even if we change the model structure from PreAct-ResNet18 to VGG-19, the conclusion coincides with our analysis. Besides, with the increase in poisoning ratio, some color patterns are closer to the anti-diagonal line, which means these defense methods can achieve better performance in this situation. Please refer to Section 4.2 in the manuscript for the analysis.

### B.3 Effect of dataset

As shown in Figure 2, we make a detailed comparison of the performance of attacks and defenses under different datasets using the PreAct-ResNet18 model and 5% poison ratio. Where the different colored bars correspond to the four datasets, the height of the bars represents the ASR, and the various subplots correspond to the multiple defenses (and no defenses). Looking down from the undefended perspective, we can see that, by and large, the effect of the attack fluctuates across the different datasets. Blended is the most stable across datasets, while BadNets has the most fluctuating effect across datasets. For BadNets, we find that CIFAR-100 and GTSRB are more complex than CIFAR-10, which leads to the decrease in effectiveness on these two datasets, but the ASR on Tiny ImageNet has rebounded significantly due to the enlargement of the trigger size. From different defense perspectives, we can find that the two methods, AC and Spectral Signature, are relatively unaffected by changes in the dataset compared with each other. In contrast, the rest of the defense methods may all have large fluctuations in their effectiveness in the face of specific attacks. Although

fluctuating, ANP has better results on CIFAR-10 for all attack methods, while ABL is also very effective on Tiny ImageNet for all attack methods.

## B.4 Effect of poisoning ratio

### B.4.1 Effect of poisoning ratio with randomness

As demonstrated in Section 4.3 in the main manuscript, in the following we will further verify the abnormal phenomenon of poisoning ratio’s effect shown in Figure 3 in the main manuscript, under the randomness of weight initialization and some methods’ mechanisms.

**Experimental setting** In the reported 8,000 pairs of evaluations, we set the random seed as 0 to fix all randomness in each evaluations, to ensure all results could be reproduced. For each evaluation plotted in Figure 3 in the main manuscript, we re-run the script with five different random seeds, and record the mean and the standard deviation of these five evaluations.

**Analysis** As shown in Figure 3, the trends of ASR curves are almost consistent with those in Figure 3 in the main manuscript, and the standard deviation (*i.e.*, the error bar) is small, indicating that the abnormal phenomenon about the poisoning ratio’s effect is not affected by the randomness. However, there are still a few special cases. For example, the error bars of some attacks under the ABL defense is very large when the poisoning ratio is low. The reason is that ABL identifies the fixed 1% of all training samples as the poisoning samples according to the training loss. However, we find that the poisoning identification accuracy is very unstable, especially when the poisoning ratio is low, leading to the large fluctuation. The standard deviations of evaluations under the NC defense are also large. As described in Section ??, NC consists of two consecutive steps, *i.e.*, firstly searching a candidate trigger to determine whether it is a backdoored model or not, then mitigating the backdoor effect through pruning. We observe that the first step is very unstable within 5 random evaluations. If the backdoored model is successfully detected, then the ASR will be reduced significantly, other keeping the high value, causing the high standard deviations of 5 random evaluations.

### B.4.2 Effect of poisoning ratio of other model architectures

In this part, we intend to give the more detailed information about the variation of ASR values against poisoning ratio on different model structures, which are VGG-19, DenseNet-161, EfficientNet-B3, and MobileNetV3-Large, respectively. The corresponding results are established in Figure 4. We have analyzed the effect of poisoning ratio in Section 4.3 in the main manuscript based on the results of Preact-ResNet18 on CIFAR-10. We found that the most ASR curves increase with the increase of poisoning ratio, while there are some curves which increase at first and then collapse dramatically. However, this phenomenon still exists for multiple model structures. It is interesting to notice that if the model structure is changed, the tendency of curves is different. The curve of NAD against BadNets can serve as an example. It keeps increasing in DenseNet-161, while increases at first and then drops down in VGG-19 and MobileNetV3-Large. Thus, it is valuable to further explore the relationship between model architecture and backdoor performance. Note that we don’t provide the results of ANP and DBD on VGG-19, as we adopt the VGG-19 architecture without the batch normalization (BN) layer (see the demonstration in Section 4.1 of the main manuscript). According to the ANP author’s comments at [https://github.com/csdongxian/ANP\\_backdoor/issues/2](https://github.com/csdongxian/ANP_backdoor/issues/2), ANP is not suitable to the model architecture without the BN layer. Besides, in our evaluations, the defense performance of DBD is not very stable on the VGG-19 without the BN layer, at its semi-supervised learning phase. Thus, we also don’t report the evaluation of DBD on VGG-19. However, we also observe that DBD performs stably on the VGG-19 model with the BN layer. The behind reason will be explored in future.

## B.5 Sensitivity to hyper-parameters

As illustrated in Section ??, we adopt a fixed setting for each attack/defense method to ensure reproducibility and fair comparison. However, the sensitivity to hyper-parameters is also a very critical metric of one method’s performance and practical usage. In the following, we pick three attack methods (*i.e.*, BadNets, SIG, InputAware), two defense methods (*i.e.*, ABL and ANP), two datasets (*i.e.*, CIFAR-10 and GTSRB), and two models (*i.e.*, PreAct-ResNet18 and VGG-19), to present a

partial analysis about the sensitivity to hyper-parameters. For each attack/defense method, we study one key hyper-parameter, such as the trigger’s patch size for BadNets, the trigger’s frequency for SIG, the mask density for Input-aware, the flooding value for ABL and the poisoning threshold for ANP.

Results are shown in Tables 3, 4, 5, and 6. For BadNets, a larger square pattern means a more vigorous attack but is also easier to find by defense methods. For SIG, higher frequency means stronger attacks and harder to defend. For Input-aware, the ASR values fluctuate a lot *w.r.t.* the mask density in the case of no defense but are relatively stable under defenses. For ABL, a higher flooding parameter often leads to worse defense performance for most attacks. For ANP, the situation is complicated. When defending BadNets and SIG, the higher threshold often leads to better defense performance, but the defense performance against Input-aware is rather stable *w.r.t.* the threshold. When comparing the performance across different model architectures and different datasets, we find that the sensitivities to hyper-parameters of each method are very diverse. Picking a good hyper-parameter for a backdoor learning method in practice is a challenge.

### B.6 Analysis of quick learning of backdoor

The quick learning phenomenon of backdoor has been observed in some previous works [9], *i.e.*, the backdoor could be quickly learned in a few epochs, for almost all backdoor attacks. However, the behind reason has not been studied. In the following, we provide a detailed analysis from the perspective of gradient. Specifically, for each epoch during the training process, we record the following information:

- Losses of training samples, clean testing samples, and poisoned testing samples;
- Accuracy on training samples, clean testing samples, and poisoned testing samples;
- Gradient signal to noise ratios (GSNR) [11] on training samples, clean train samples, and poisoned training samples averaged over model parameters;
- Norms of average gradient on total training samples, clean training samples, and poisoned training samples;
- Pairwise cosine similarities between average gradients on total training samples, clean training samples, and poisoned training samples.

As shown in Figure 5, we report the results of 5 backdoor attacks, including BadNets, Blended, SSBA, LC and LF with poisoning ratio 10%, on the CIFAR-10 dataset and the PreAct-ResNet18 model. As shown in the first column, the testing loss of poisoned samples drops quickly in the early stages of training and converges to a low value, while the testing loss of clean samples drops at a slower rate and converges to a much larger value. It verifies the quick learning phenomenon of backdoor under these five backdoor attack methods. As shown in the third column, we first observe that the GSNR of poisoned samples is significantly larger than the GSNR of clean samples at the early stages. The high GSNR values of poisoned samples indicate that the backdoor has better generalization performance and is consistent with the higher accuracy (ASR) and lower loss on poisoned testing samples. Secondly, we notice that the norm of the gradient on poisoned samples is much larger than the norm of the gradient on clean samples in early epochs, as shown in the fourth column. Consequently, the cosine similarity between gradients on total training samples and poisoned training samples is significantly larger than the cosine similarity between gradients on clean training samples and poisoned training samples, though the number of poisoned samples is much smaller than the number of clean samples.

### B.7 Analysis of backdoor forgetting

From the above analysis about the quick learning of backdoor in Section B.6, we get one impression that the backdoored model memorizes the poisoned samples quickly and stably. To obtain more insights about the inner mechanism of backdoor learning, we adopt the concept of *forgetting event* [23] to characterize the learning dynamics during the training process. Specifically, one forgetting event is recorded when a correctly predicted training sample at the current epoch is incorrectly predicted at the next epoch. Formally, given a training sample  $(\mathbf{x}, y)$ , where  $\mathbf{x}$  is input feature,  $y$  is ground-truth label. If  $\mathbf{x}$  is correctly predicted at the epoch  $t$ , *i.e.*,  $f_{\theta^t}(\mathbf{x}) = y$ , but is misclassified at epoch  $t + 1$ , *i.e.*,  $f_{\theta^{t+1}}(\mathbf{x}) \neq y$ , where  $f$  denotes the model and  $\theta^t, \theta^{t+1}$  are model parameters at the epoch  $t$  and  $t + 1$ , then a forgetting event is recorded for this sample.

Specifically, we count the number of forgetting events for clean and poisoned training samples, respectively, on CIFAR-10 with Preact-ResNet18 backbone. The distributions of forgetting events of clean and poisoned samples are shown in Figure 6. The results show that:

- The forgetting events of clean training samples follow an exponential distribution, and are similar among different cases.
- For poisoned training samples: **1)** when the poisoning ratio is low (*e.g.*, 0.1%, 0.5%), the forgetting numbers of poisoned samples are often larger than those of clean samples; **2)** when the poisoning ratio is high (*e.g.*, 5%, 10%), the forgetting numbers of poisoned samples are often smaller than those of clean samples.

The above observations are compatible with our high-level observation that the backdoor attack with higher poisoning ratios could quickly learn the stable mapping from the poisoned samples to the target class. Moreover, the forgetting event provides a fine-grained tool to analyze the contribution of each individual training sample, which could facilitate the development of more advanced backdoor attack and defense methods.

### B.8 Analysis of trigger generalization of backdoor attacks

In all existing backdoor attacks, there is a default assumption that the triggers used in both backdoor training and backdoor testing are exactly same. However, we find that an interesting property in backdoor learning that the backdoored model trained with one trigger could be also activated by other triggers. We name it as **trigger generalization**.

In this following, we take the Blended attack as an example to study the trigger generalization. Specifically, we set the trigger transparency to different values in training and testing phase, including 10%, 20%, and 30%. As shown in the figure 7, we obtain the following observations: (1) If a more obvious trigger (*i.e.*, high transparency 30%) is applied during training, then the backdoor will not be easily activated by a different trigger with lower transparency; (2) If using a less obvious backdoor triggers (*i.e.*, low transparency 10%) in the training phase, then the backdoor can be successfully activated by the trigger with higher transparency (*e.g.*, 20% or 30%) in the test phase.

In addition to the above example, we find that the trigger generalization is a common property of the backdoor models under several backdoor attacks. For example, the trigger in SSBA [10] is a string, and its backdoored model could not only be activated by its training trigger. However, we find that the model is likely to be activated by many other strings with the same length. Exploring the behind reason of trigger generalization is important for us to better understand the backdoor mechanism. Moreover, we notice that there have been some attempts to utilize the backdoor as the technique to protect the intellectual property (IP) of AI models or datasets, based on the unique mapping from the trigger to the target class. However, due to the trigger generalization, the uniqueness of the training trigger no longer exists, which undermines the legitimacy backdoor learning in IP protection. Thus, the study of trigger generation is also important for the usage of backdoor learning in practice. We will provide more analysis about trigger generalization in BackdoorBench in future.

### B.9 Evaluation on vision transformer

Until now, the reported 8,000 pairs of evaluations in BackdoorBench are all conducted on the convolutional neural networks. In the following, we expand the evaluations to another popular family of models, *i.e.*, vision transformer (ViT), which has shown superior performance on many vision tasks (*i.e.*, image classification, object detection, semantic segmentation). In our evaluations, a initial checkpoint of ViT that is pre-trained on ImageNet [4] is downloaded from <https://github.com/pytorch/vision/tree/main/references/classification>. We then fine-tune this pre-trained checkpoint on the poisoned CIFAR-10 dataset. Note that the input size of ViT is  $224 \times 224$ , while the size of raw images in CIFAR-10 is  $32 \times 32$ . Thus, in data poisoning based attack, we firstly insert the trigger into the raw image, then re-scale the poisoned image to the size  $224 \times 224$ .

We evaluate ViT-Base model with  $16 \times 16$  input patch size (ViT-b-16) on CIFAR-10 with 10% poison ratio, where the settings of all hyper-parameters are same with those for learning other models, as shown in Tables 1 and 2. The backdoor evaluation results are summarized in Table 7. As a baseline, the accuracy of fine-tuning ViT-b-16 on the clean dataset with same hyper-parameters is 96.56%. According to Table 7 and the comparison with the evaluations on other models, we have the following



observations. **1)** In the case of no defense, the ASR of ViT-b-16 is still very high under all evaluated backdoor attacks, revealing that the ViT model architecture is also vulnerable to backdoor attacks. **2)** The evaluated three defense methods show very poor performance for the attack on ViT-b-16. In terms of FT and NC, although the ASR is reduced significantly, the clean accuracy is also downgraded. In terms of ABL, the ASR doesn't decrease for most attacks, with the only exception for BadNets of which the model after defense is fully degenerated. It implies that the effective defense that have been verified on the CNN architecture may not suitable for the ViT architecture. It inspires us to develop more effective defense methods for the ViT architecture specially. More evaluations and analysis about the ViT architecture will be added in BackdoorBench in future.

## B.10 Evaluation on ImageNet

Due to the high computational and memory costs, one of the benchmark datasets of image classification, *i.e.*, ImageNet [4] with 1,000 classes, has rarely been evaluated in existing backdoor learning works. We plan to provide comprehensive evaluations of backdoor learning methods on ImageNet, to find whether there are some unique challenges for backdoor learning on large-scale datasets.

Here, we provide some partial evaluations, including BadNets and Blended attack with 0.1% poison ratio on ImageNet and the PreAct-ResNet18 model, as shown in Table 8. Note that due to the 1,000 classes, we do not set a higher poisoning ratio to ensure that the number of poisoned samples is not much larger than the number of clean samples of the target class. Both Badnets and Blended show good attack performance with high ASR and C-Acc. Compared with the baseline model, *i.e.*, PreAct-ResNet18 trained on the clean ImageNet dataset (please refer to [https://paperswithcode.com/sota/image-classification-on-imagenet?tag\\_filter=3](https://paperswithcode.com/sota/image-classification-on-imagenet?tag_filter=3)), there is a slight drop of C-Acc, from 72.33% to 69.22%. More backdoor attack and defense evaluations on ImageNet will be added to our BackdoorBench in the future.

## B.11 Visualization

### B.11.1 Individual visualization tools

Here we provide three visualization tools to analyze each individual image.

**Gradient-weighted class activation mapping (Grad-CAM)** [19] explains the contribution of each pixel to the prediction of one image, based on the gradient of the logit of one class *w.r.t.* each pixel. Note that in the codebase of BackdoorBench, we implement a variant of Grad-CAM, called FullGrad [22].

**Shapley Value** [13] is another popular interpretation tool that assigns an importance factor to each pixel for a particular prediction. Inspired by the cooperative game theory, the competition among pixels is also taken into account in the computation of each individual importance factor.

**Frequency saliency map (FSM)** is our innovative visualization method for viewing the contribution of every Fourier basis to model classification. Consider an image classification task with  $S$  classes. Let  $\mathbf{x}$  be a clean image with size  $H \times W \times C$  and  $\tilde{\mathbf{x}} = \mathfrak{F}(\mathbf{x})$  be the corresponding frequency spectrum with  $\mathfrak{F}$  being the channel-wise Discrete Fourier Transform (DFT) operator. Let  $\mathbb{C}$  be the set of complex numbers. Denote the classifier by  $f : \mathbb{R}^{H \times W \times C} \rightarrow \mathbb{R}^S$ . We define  $F : \mathbb{C}^{H \times W \times C} \rightarrow \mathbb{R}^S$  as the corresponding classifier in the frequency domain, which means

$$F(\tilde{\mathbf{x}}) = f(\mathbf{x}) = f(\mathfrak{F}^{-1}(\tilde{\mathbf{x}})), \quad (1)$$

where  $\mathfrak{F}^{-1}$  is the channel-wise Inverse Discrete Fourier Transform (IDFT) operator, which means  $\mathfrak{F}^{-1}(\tilde{\mathbf{x}})(u, v) = \sum_{h=0}^{H-1} \sum_{w=0}^{W-1} \tilde{\mathbf{x}}(h, w) e^{2\pi i(\frac{uh}{H} + \frac{vw}{W})}$ .

Inspired by the saliency map in the RGB space, we intend to establish the connection between model prediction and image's frequency spectrum by the gradient. According to the chain's rule, we can estimate the gradient *w.r.t.* the frequency spectrum as follows

$$\begin{aligned} \frac{\partial F_s(\tilde{\mathbf{x}})}{\partial \tilde{\mathbf{x}}(u, v, c)} &= \sum_{h=0}^{H-1} \sum_{w=0}^{W-1} \sum_{c'=0}^{C-1} \frac{\partial f_s(\mathbf{x})}{\partial \mathbf{x}(h, w, c')} \cdot \frac{\partial \mathbf{x}(h, w, c')}{\partial \tilde{\mathbf{x}}(u, v, c)} \\ &= \sum_{h=0}^{H-1} \sum_{w=0}^{W-1} \frac{\partial f_s(\mathbf{x})}{\partial \mathbf{x}(h, w, c)} e^{2\pi i(\frac{uh}{H} + \frac{vw}{W})}, \end{aligned} \quad (2)$$

where  $f_s$  means the logit output of the model  $f$  w.r.t. the  $s$ -th class.

The obtained gradients are then averaged over channels and normalized to form the FSM, as follows

$$FSM^*(u, v) = \frac{1}{C} \sum_{c=0}^{C-1} \frac{\partial F_s(\tilde{\mathbf{x}})}{\partial \tilde{\mathbf{x}}(u, v, c)}, \quad (3)$$

$$FSM(u, v) = \frac{FSM^*(u, v) - FSM_{\min}^*}{FSM_{\max}^* - FSM_{\min}^*}, \quad (4)$$

where we use min-max normalization to normalize the FSM image.

### B.11.2 Visualization results

In the following, we present some visualization results using the above three tools to understand the inner mechanism of backdoor learning better. Specifically, we train the PreAct-ResNet18 model under various backdoor attacks and defenses with the poisoning ratio 5%, on 3 datasets, including CIFAR-100, GTSRB, and Tiny ImageNet. Then, we randomly select a poisoned sample from the test set and show its visualization. The visualization results using Shapley Value and Grad-CAM are shown in Figures 8 to 13, respectively.

The frequency saliency map (FSM) visualization results are shown in Figure 14, where low-frequency components are shifted into the central regions. In contrast, high-frequency components are distributed in surrounding regions. The first column displays the studied poisoned images generated by different attack methods, including BadNet, Blended, SSBA, WaNet, and LF. The second column shows the contribution of each frequency basis to the backdoored model’s prediction. It tells that most backdoor models pay attention to high-frequency regions, while the model under the LF attack makes the model concentrates more on low-frequency regions. Besides, we can see an apparent difference between SSBA and WaNet in the frequency domain, even though their spatial images look similar. These observations demonstrate the potential usage of FSM in trigger or backdoor detection, which will be explored in our future work. The remaining columns show the contribution of each frequency basis under various defense methods. For example, the contribution of some high-frequency regions is enormous for the poisoned image in a BadNet model without defense. However, after conducting FT on this backdoor model, the low-frequency regions regain attention from the model. It explains well that FT can effectively remove backdoors embedded by BadNet and WaNet (see Figure 3 in the main manuscript). We plan to explore more backdoor learning properties from the frequency domain perspective.

## C BackdoorBench in Natural Language Processing

Apart from the analysis of backdoor attack and defense methods in computer vision, we also expand our benchmark to the field of Natural Language Processing (see [https://github.com/SCLBD/BackdoorBench/tree/main/backdoorbench\\_nlp](https://github.com/SCLBD/BackdoorBench/tree/main/backdoorbench_nlp)). We implement two state-of-the-art backdoor attack methods (*i.e.*, LWS [18] and HiddenKiller [17]) and one defense method (*i.e.*, Onion [16]) in NLP as a complement to the original BackdoorBench. We closely follow the original implementation of the attack and defense methods and make necessary changes to unify all the methods in our benchmark. We choose BERT [5] as the model to be poisoned. All the experiments are conducted on three widely-used datasets for text classification tasks, including Stanford Sentiment Treebank (SST-2) [21], Offensive Language Identification Dataset (OLID) [28] and AG’s News [29]. For all experiments, the poison rate is set to be 5% and the default target label is 1. For LWS, the bar for ONION for each dataset is set to be the recommended value in the original implementation. All results are reported in Table 9.

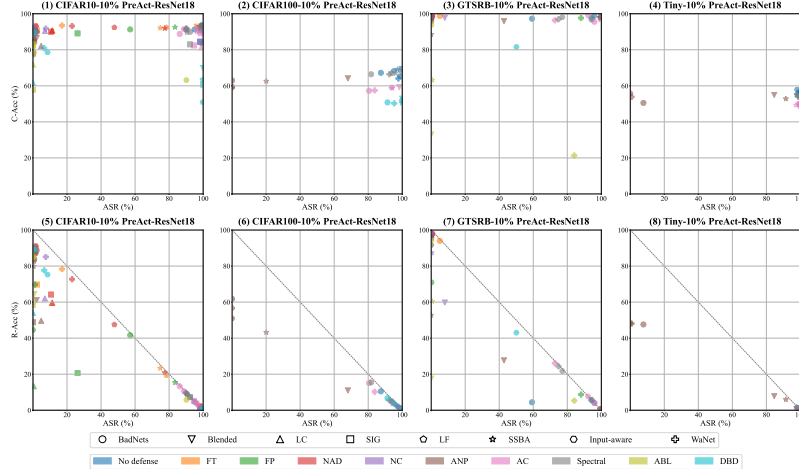
We can find that the two chosen attack methods can both achieve high attack success rate even at a low poisoning ratio. However, the defense performance of ONION against two SOTA attack methods is not quite satisfactory. The possible reason is that ONION aims to find out obvious outliers in each sentence, but both HiddenKiller and LWS are invisible methods which do not rely on special tokens as triggers. In the future, we will also keep updating latest backdoor attack and defense methods in the NLP field into our benchmark.

## **D Reproducibility**

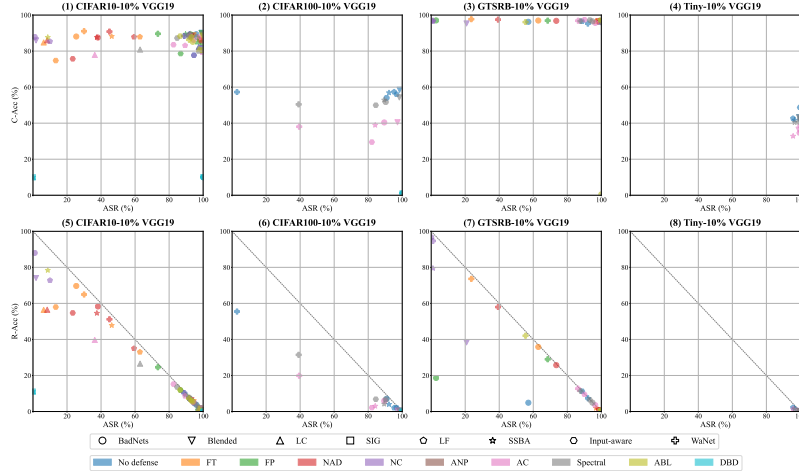
All evaluation results in BackdoorBench can be easily reproduced, just running the scripts provided in the github repository <https://github.com/SCLBD/BackdoorBench>, with the hyper-parameter settings presented in Tables 1 and 2. All evaluated datasets and model architectures are publicly and freely available. Besides, we also compress all codes into one file as a part of the supplementary materials.

## **E License**

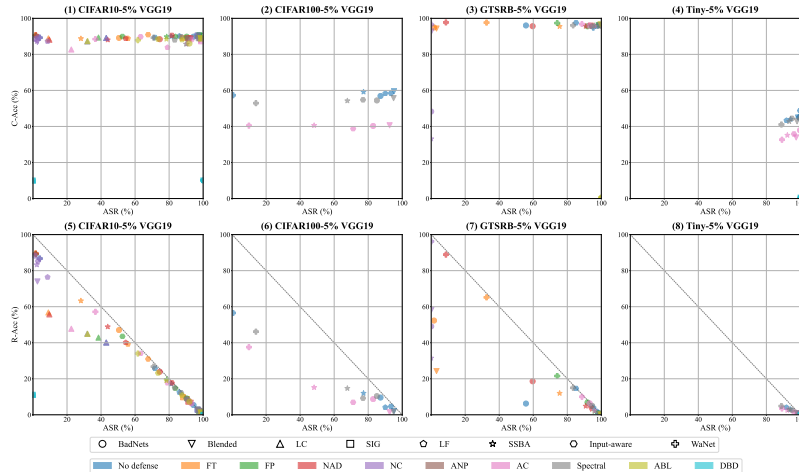
This repository is licensed by The Chinese University of Hong Kong, Shenzhen and Shenzhen Research Institute of Big Data under Creative Commons Attribution-NonCommercial 4.0 International Public License (identified as CC BY-NC-4.0 in SPDX, see <https://spdx.org/licenses/>). More details about the license could be found in <https://github.com/SCLBD/BackdoorBench/blob/main/LICENSE>.



(a) Attack-defense pairs with PreAct-ResNet18 and 10% poisoning ratio on CIFAR-10.



(b) Attack-defense pairs with VGG-19 and 10% poisoning ratio on CIFAR-10.



(c) Attack-defense pairs with VGG-19 and 5% poisoning ratio on CIFAR-10.

Figure 1: Performance distribution of attack-defense pairs on different model structure and poisoning ratios. A successful attack method should be high C-Acc and ASR, while a successful defense method should be high C-Acc and low ASR. Besides, if the reduced ASR value equals to the increased R-Acc, the color patters would be close to the anti-diagonal line.

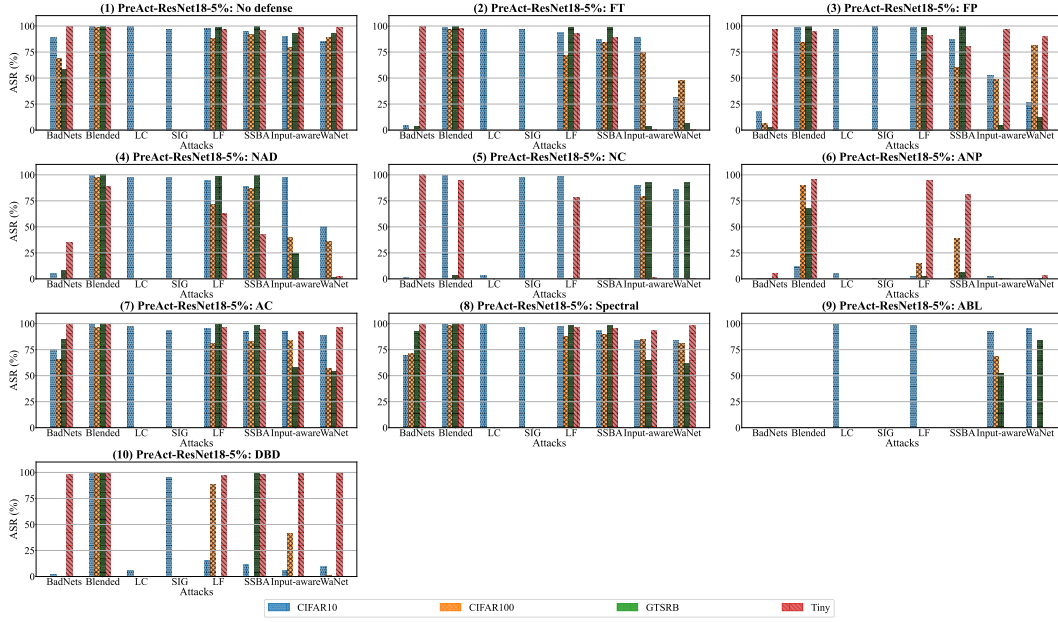


Figure 2: The effects of different datasets on backdoor learning. Note that for the clean-label (*i.e.*, LC [20] and SIG [1]) attack, the number of poisoned samples must be less than the target class size, thus it may be not applied to the case of high poisoning ratios.)

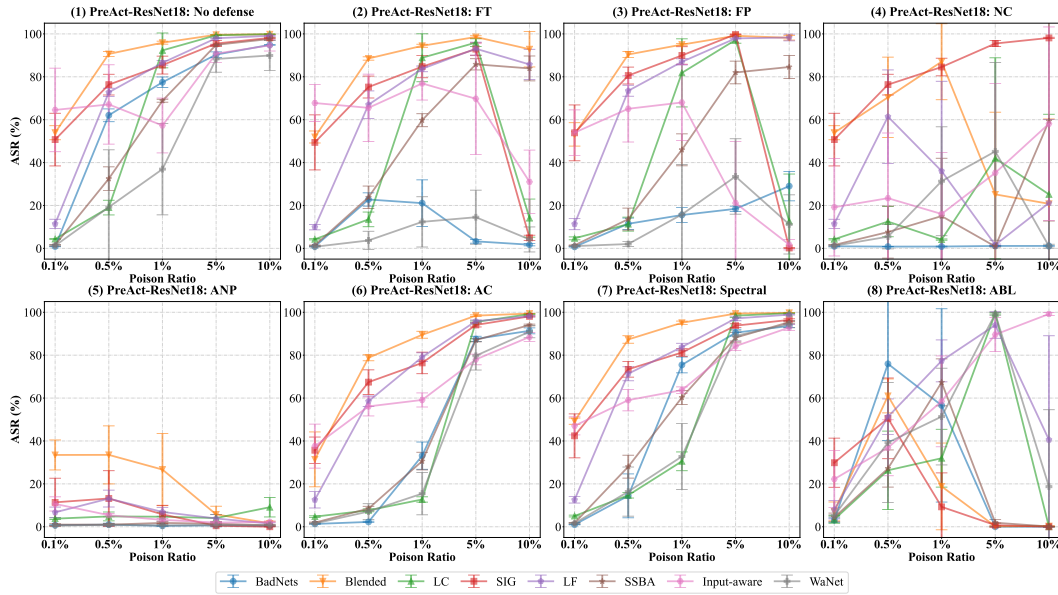
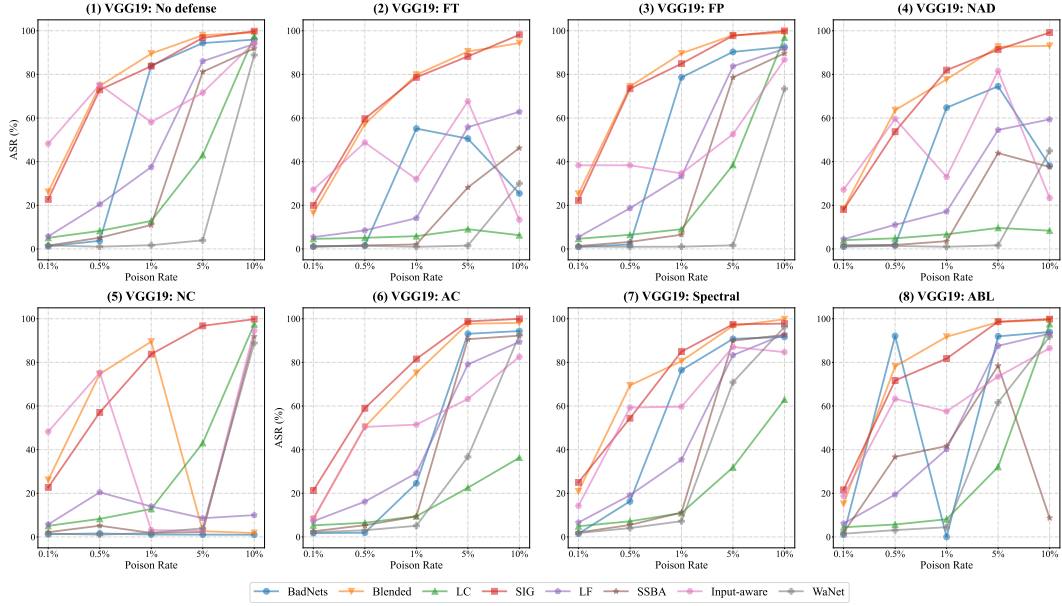
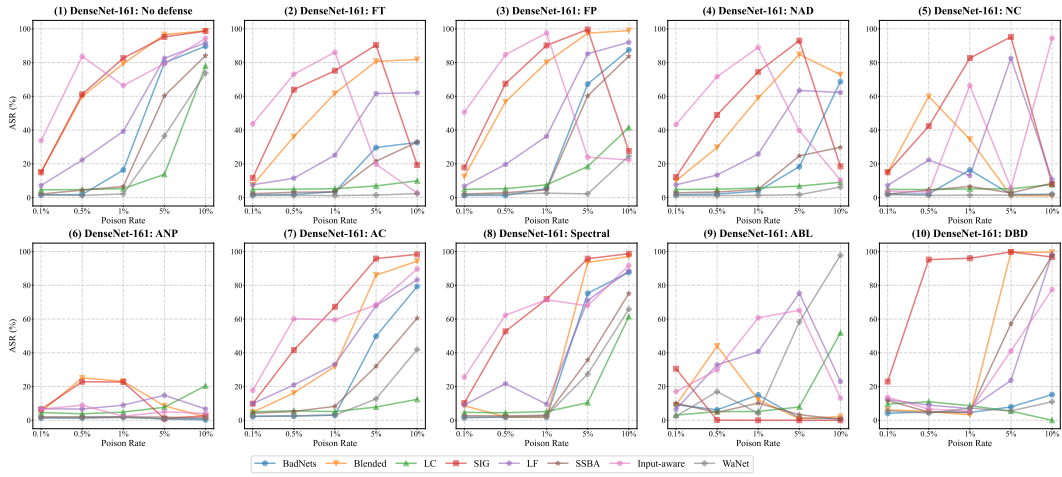


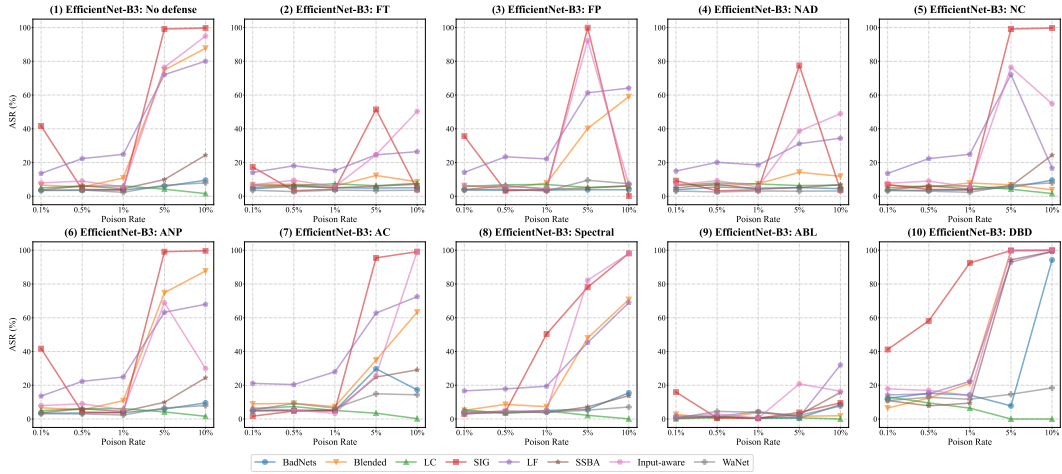
Figure 3: The effects of different poisoning ratios of backdoor learning with 5 random seeds.



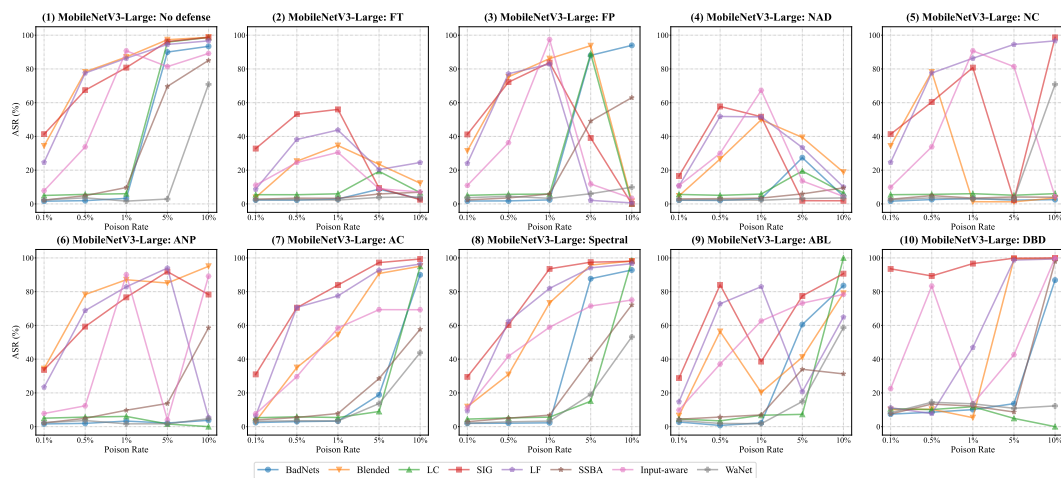
(a) The variation of ASR on different poisoning ratios with VGG-19 and CIFAR-10.



(b) The variation of ASR on different poisoning ratios with DenseNet-161 and CIFAR-10.



(c) The variation of ASR on different poisoning ratios with EfficientNet-B3 and CIFAR-10.



(d) The variation of ASR on different poisoning ratios with MobileNetV3-Large and CIFAR-10.

Figure 4: The effect of different ratios on backdoor learning. From (a) to (d), the structure of models are different. In the condition of no defense, the higher poisoning ratio, the higher ASR value. In the defense situation, some ASR curves raise with the increase of poisoning ratio, while some curves go up first and then sharply drop down. It could also be noticed that the performances of same defense method on different model structures are distinctive, *i.e.*, ABL on VGG-19 and DenseNet-161. Note that we don't provide the results of ANP and DBD on VGG-19, and the reason is illustrated in Section B.4.2.





Table 4: Sensitivity results on CIFAR10 with VGG-19.

	No Defense		No Defense		ABL		ABL		ABL		ABL		ABL		ABL		ABL			
	C-Acc (%)	ASR (%)	R-Acc (%)	R-Acc (%)	C-Acc (%)	ASR (%)	R-Acc (%)	R-Acc (%)	C-Acc (%)	ASR (%)	R-Acc (%)	R-Acc (%)	C-Acc (%)	ASR (%)	R-Acc (%)	R-Acc (%)	C-Acc (%)	ASR (%)	R-Acc (%)	
BadNets 1	85.79	76.44	22.46	26.70	84.14	70.67	26.70	43.34	53.99	45.03	36.38	10.00	100.00	0.00	10.00	100.00	0.00	10.00	100.00	0.00
BadNets 2	87.26	93.76	5.80	23.83	84.55	73.70	23.83	9.62	83.58	88.21	10.67	78.82	85.81	12.76	10.00	100.00	0.00	10.00	100.00	0.00
BadNets 3	88.75	95.77	3.93	7.87	85.93	91.29	7.87	8.56	84.31	93.03	6.28	83.75	92.74	6.54	79.32	91.17	7.73	79.32	91.17	7.73
BadNets 4	89.50	96.73	3.08	5.52	86.41	93.91	5.52	6.80	87.24	94.44	5.02	85.75	95.06	4.50	85.69	94.64	4.90	85.69	94.64	4.90
BadNets 5	89.87	97.03	2.78	3.53	88.04	96.18	3.53	3.02	87.79	96.98	2.76	88.22	97.04	2.73	88.28	97.50	2.27	88.28	97.50	2.27
SIG 1	82.20	41.61	53.47	41.06	80.54	41.06	51.78	49.36	80.68	44.10	50.00	80.33	44.80	49.26	80.14	45.06	47.74	80.14	45.06	47.74
SIG 2	81.92	74.87	21.84	23.38	80.16	72.33	23.38	20.68	80.39	78.21	19.44	80.02	82.79	14.37	80.05	83.04	14.23	80.05	83.04	14.23
SIG 3	81.44	88.27	8.38	10.23	80.25	87.62	10.23	9.21	80.14	87.14	10.50	80.46	85.22	11.08	80.40	90.36	7.99	80.40	90.36	7.99
SIG 4	81.44	98.03	1.42	1.32	80.73	98.43	1.32	1.40	80.45	99.27	0.61	80.55	98.76	1.01	78.89	94.52	4.13	78.89	94.52	4.13
SIG 5	82.16	99.50	0.30	0.42	80.04	99.34	0.42	0.63	80.17	99.17	0.72	80.58	98.56	1.08	80.01	98.41	1.24	80.01	98.41	1.24
Input-aware 0.032	88.48	96.77	2.89	8.19	88.27	90.49	8.19	6.90	24.65	100.00	0.00	86.56	7.87	73.64	10.00	100.00	0.00	10.00	100.00	0.00
Input-aware 0.048	88.86	94.98	4.44	16.21	88.31	81.49	16.21	17.93	88.78	72.57	24.50	10.00	100.00	0.00	10.00	100.00	0.00	10.00	100.00	0.00
Input-aware 0.064	88.71	85.72	12.91	13.64	88.58	84.52	13.64	1.03	10.00	100.00	0.00	50.51	0.00	37.17	10.00	100.00	0.00	10.00	100.00	0.00
Input-aware 0.08	88.54	97.00	2.61	13.40	87.81	84.16	13.40	11.50	88.04	88.20	10.16	87.66	87.62	10.60	10.00	100.00	0.00	10.00	100.00	0.00
Input-aware 0.096	89.11	96.52	3.36	15.76	88.01	81.91	15.76	14.17	88.38	84.06	14.08	87.39	6.69	73.10	10.00	100.00	0.00	10.00	100.00	0.00



Table 6: Sensitivity results on GTSRB with VGG-19.

	No Defense	No Defense	No Defense	ABL	ABL	ABL	ABL	ABL	ABL	ABL	ABL	ABL	ABL	ABL	ABL	ABL	ABL	ABL
	ACC	ASR	RA	0.1	0.1	0.1	0.5	0.5	0.5	0.9	0.9	0.9	1.3	1.3	1.3	1.7	1.7	1.7
	ACC	ASR	RA	ACC	ASR	RA	ACC	ASR	RA	ACC	ASR	RA	ACC	ASR	RA	ACC	ASR	RA
BadNets 1	96.17	89.01	9.84	96.13	91.85	7.81	0.48	100.00	0.00	0.48	100.00	0.00	0.48	100.00	0.00	0.48	100.00	0.00
BadNets 2	95.20	93.50	5.91	0.48	100.00	0.00	0.48	100.00	0.00	0.48	100.00	0.00	0.48	100.00	0.00	57.59	96.01	2.79
BadNets 3	96.14	94.52	5.18	0.48	100.00	0.00	70.24	42.78	41.52	0.48	100.00	0.00	0.48	100.00	0.00	0.48	100.00	0.00
BadNets 4	94.35	96.38	3.31	96.26	95.55	4.27	96.26	93.09	6.57	95.56	93.35	6.38	92.91	93.20	6.26	90.07	91.03	8.23
BadNets 5	95.77	96.40	3.31	96.63	96.36	3.37	96.63	96.86	2.80	96.31	96.25	3.45	95.93	96.13	3.64	96.38	96.12	3.75
Input-aware 0.032	97.51	88.97	10.61	96.83	53.95	43.74	96.63	58.14	40.04	95.94	59.78	38.02	0.48	100.00	0.00	0.48	100.00	0.00
Input-aware 0.048	97.44	60.09	39.36	96.35	61.02	37.30	8.99	99.83	0.17	0.48	100.00	0.00	0.48	100.00	0.00	0.48	100.00	0.00
Input-aware 0.064	96.55	89.63	9.99	96.67	54.98	43.28	96.71	51.89	45.87	0.48	100.00	0.00	0.48	100.00	0.00	0.48	100.00	0.00
Input-aware 0.08	96.89	86.45	13.43	96.46	51.20	46.48	96.68	44.12	53.95	0.48	100.00	0.00	0.48	100.00	0.00	0.48	100.00	0.00
Input-aware 0.096	96.25	67.77	31.54	96.16	49.47	47.90	0.48	100.00	0.00	0.48	100.00	0.00	0.48	100.00	0.00	0.48	100.00	0.00

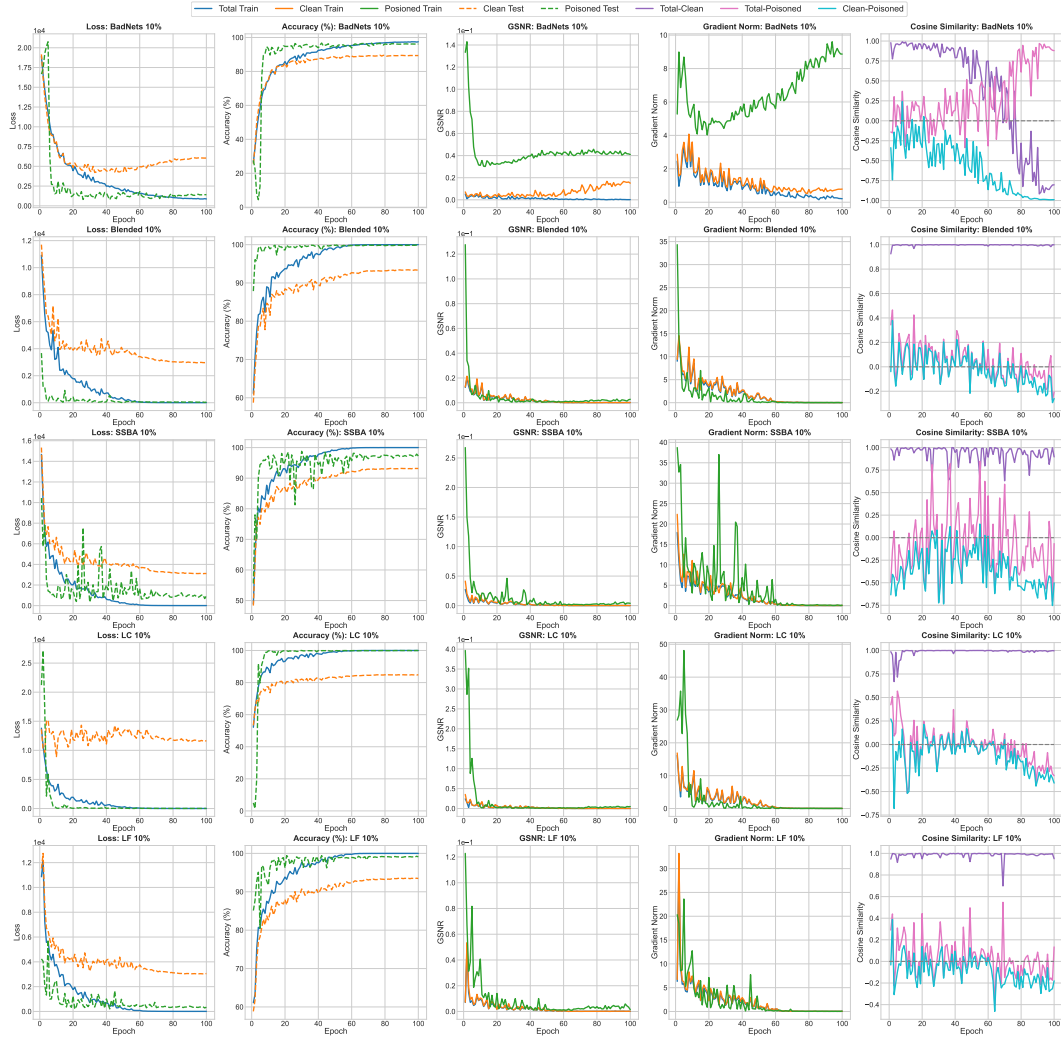


Figure 5: Analysis of the quick learning phenomenon based on the details of the backdoor learning process. From the first column to the last column, we report the curves of loss, accuracy, GSNR, gradient norm and cosine similarity of gradients in poisoned samples, clean samples and all samples, respectively.

Table 7: Evaluation of ViT-b-16 on CIFAR-10 with 10% poisoning ratio.

Backdoor Attack	No defense			FT			NC			ABL		
	C-Acc (%)	ASR (%)	R-Acc (%)	C-Acc (%)	ASR (%)	R-Acc (%)	C-Acc (%)	ASR (%)	R-Acc (%)	C-Acc (%)	ASR (%)	R-Acc (%)
BadNets	94.58	94.11	5.66	42.00	8.81	38.60	36.09	11.12	33.60	10.00	0.00	11.11
Blended	96.47	99.72	0.28	41.36	7.24	36.29	44.92	3.90	39.40	96.88	99.81	0.19
LC	86.87	99.84	0.16	33.03	16.40	24.36	40.68	11.17	30.74	87.33	99.76	0.23
SIG	87.01	92.60	7.28	44.46	0.78	25.92	87.01	92.60	7.28	87.60	86.47	13.39
SSBA	96.30	97.58	2.34	45.35	7.62	43.49	46.24	8.17	43.60	96.78	98.22	1.70
Input-aware	92.10	96.28	3.51	89.56	25.14	69.91	43.13	7.79	37.81	96.70	93.21	6.58

Table 8: BadNets and Blended attack result on ImageNet

	C-Acc	ASR	R-Acc
BadNets	69.21	75.86	0.33
Blended	69.24	98.59	0.11

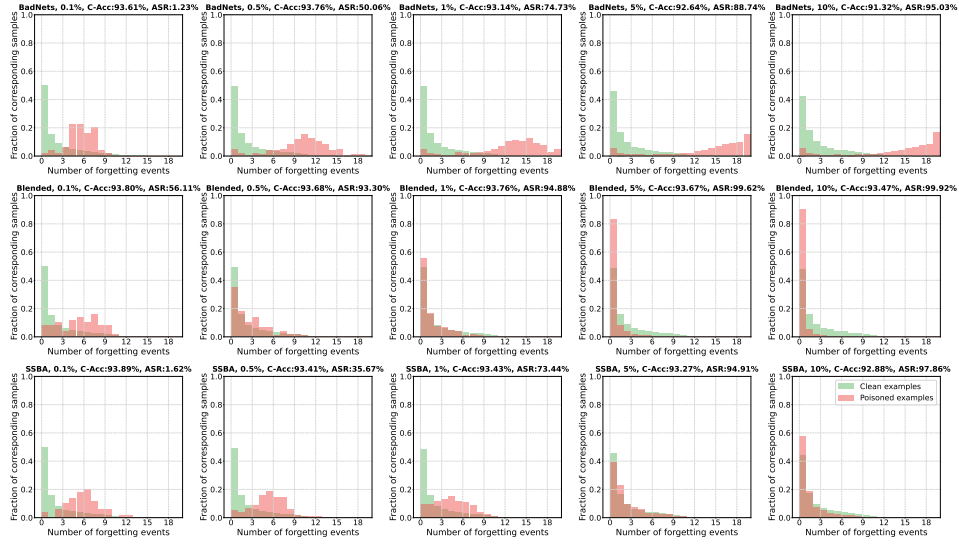


Figure 6: Distributions of forgetting events of clean training examples and poisoned examples in CIFAR-10 with Preact-Resnet18 backbone.



Figure 7: Analysis of trigger generalization in Blended attack: a) the training trigger with 10% transparency; b) the training trigger with 20% transparency; c) the training trigger with 30% transparency. For each case, we evaluate the attack success rate of testing triggers with the transparency 10%, 20%, and 30%, respectively.

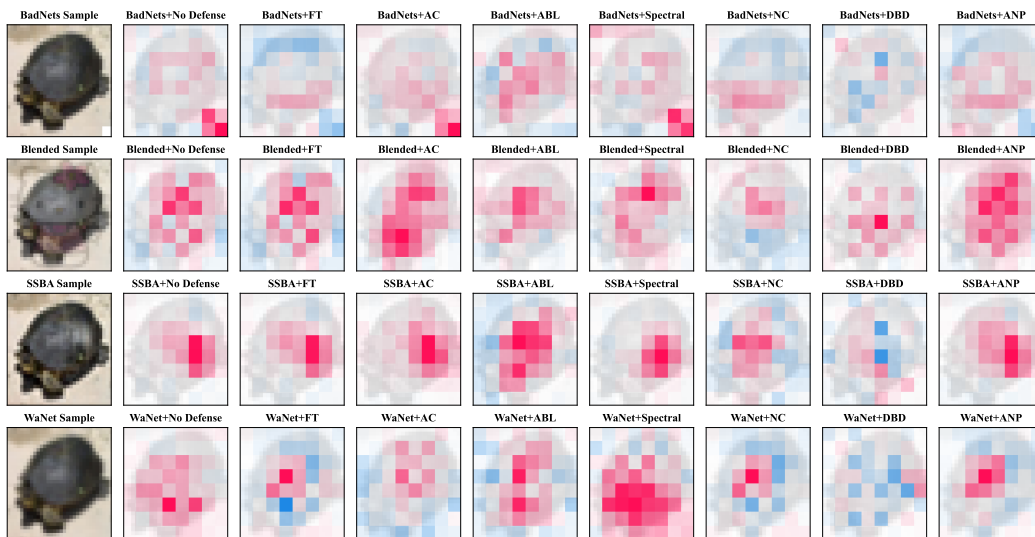


Figure 8: Shapley Value visualization of regions contributed to model decision under different attack methods and defense methods with PreAct-ResNet18 and 5% poisoning rate on CIFAR-100.

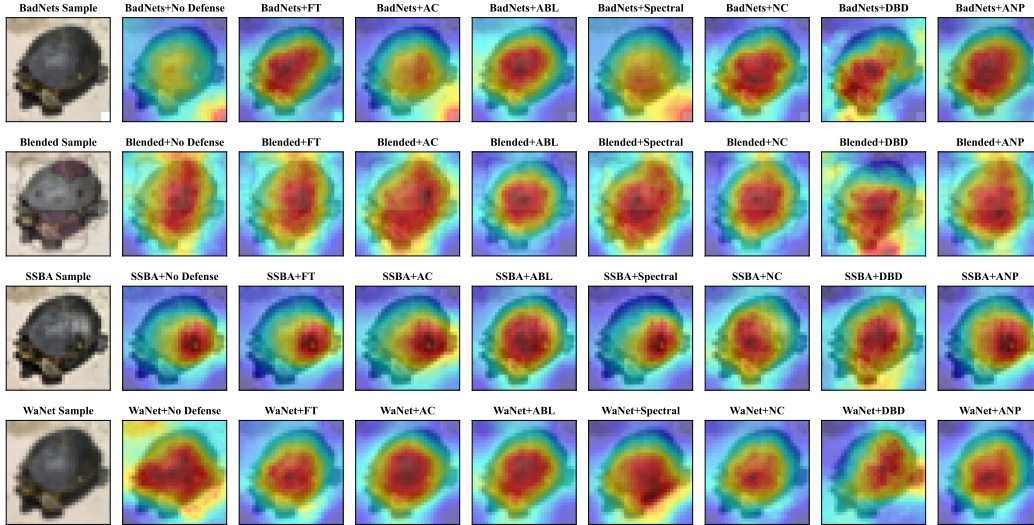


Figure 9: Grad-CAM visualization of regions contributed to model decision under different attack methods and defense methods with PreAct-ResNet18 and 5% poisoning rate on CIFAR-100.

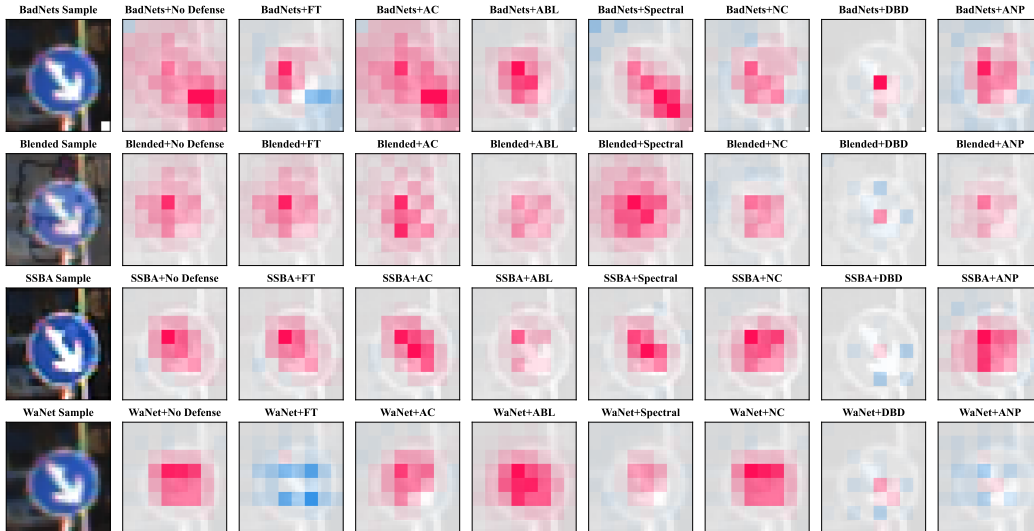


Figure 10: Shapley Value visualization of regions contributed to model decision under different attack methods and defense methods with PreAct-ResNet18 and 5% poisoning rate on GTSRB.

Table 9: Full results on backdoor attack and defense methods for BERT

		No defense			Onion		
Backdoor Attack	Dataset	C-Acc (%)	ASR (%)	R-Acc (%)	C-Acc (%)	ASR (%)	R-Acc (%)
LWS	SST-2	89.01	94.08	4.28	86.20	90.42	9.58
LWS	OLID	82.67	97.92	0.83	79.07	96.77	3.23
LWS	AgNews	93.11	99.19	0.61	92.10	68.03	10.97
HiddenKiller	SST-2	90.34	88.93	11.07	85.67	88.27	11.73
HiddenKiller	OLID	82.19	97.42	2.59	81.37	96.12	3.88
HiddenKiller	AgNews	93.49	98.67	1.12	92.05	95.16	4.21

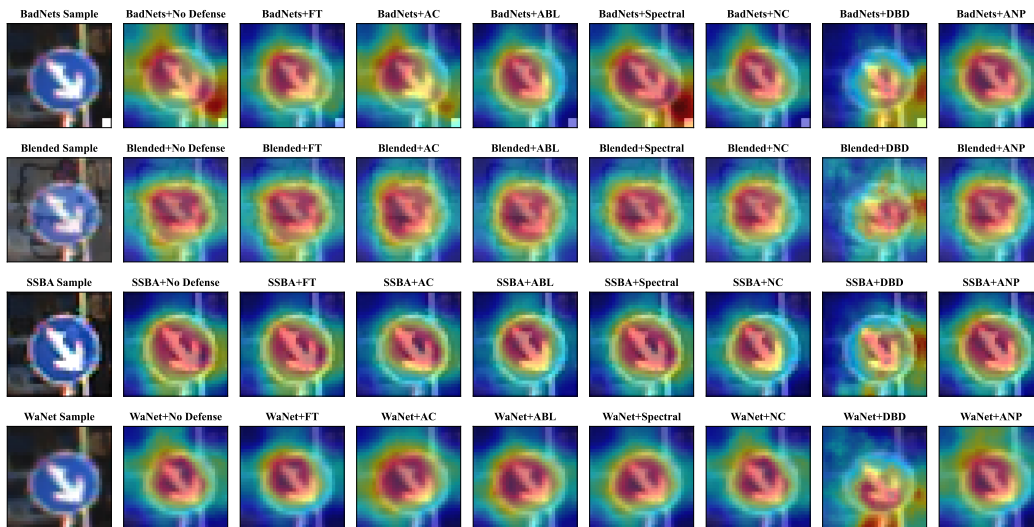


Figure 11: Grad-CAM visualization of regions contributed to model decision under different attack methods and defense methods with PreAct-ResNet18 and 5% poisoning rate on GTSRB.

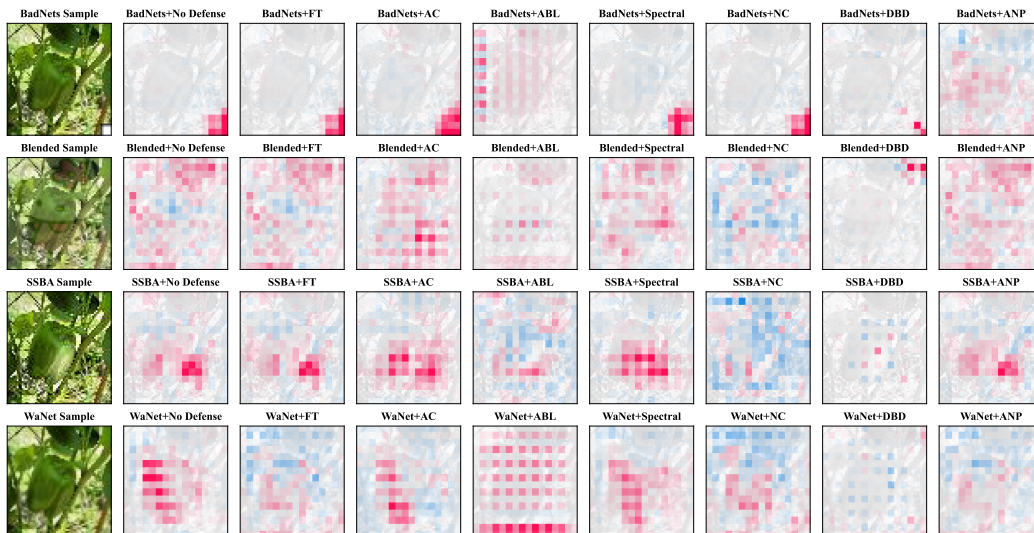


Figure 12: Shapley Value visualization of regions contributed to model decision under different attack methods and defense methods with PreAct-ResNet18 and 5% poisoning rate on Tiny ImageNet.

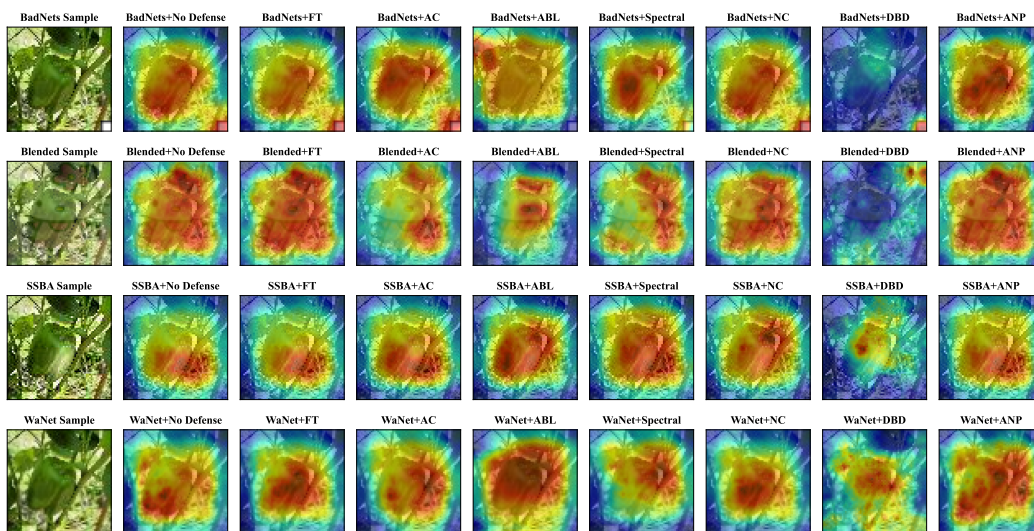


Figure 13: Grad-CAM visualization of regions contributed to model decision under different attack methods and defense methods with PreAct-ResNet18 and 5% poisoning rate on Tiny ImageNet.

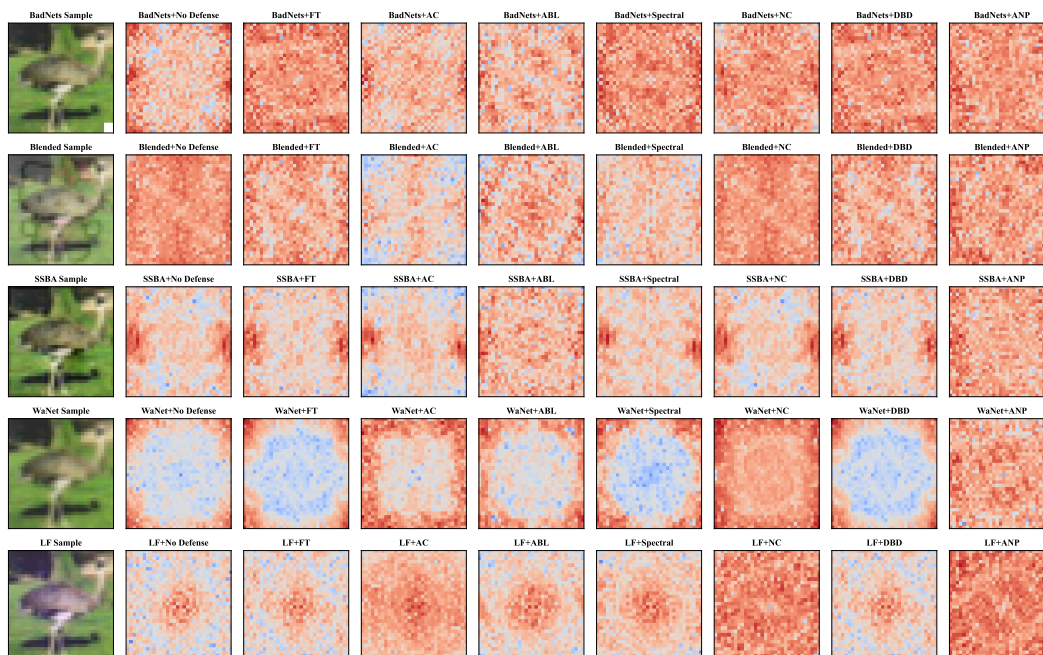


Figure 14: Frequency visualization of regions which contribute to model classification under different attack methods and defense methods with PreAct-ResNet18 and 10% poisoning rate on CIFAR-10. Warmer color was allocated to higher value whose range is 0 to 255.













## References

- [1] Mauro Barni, Kassem Kallas, and Benedetta Tondi. A new backdoor attack in cnns by training set corruption without label poisoning. In *2019 IEEE International Conference on Image Processing*, 2019.
- [2] Bryant Chen, Wilka Carvalho, Nathalie Baracaldo, Heiko Ludwig, Benjamin Edwards, Taesung Lee, Ian Molloy, and Biplav Srivastava. Detecting backdoor attacks on deep neural networks by activation clustering. In *The AAAI Conference on Artificial Intelligence Workshop*, 2019.
- [3] Xinyun Chen, Chang Liu, Bo Li, Kimberly Lu, and Dawn Song. Targeted backdoor attacks on deep learning systems using data poisoning. *arXiv preprint arXiv:1712.05526*, 2017.
- [4] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [5] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2019.
- [6] Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Identifying vulnerabilities in the machine learning model supply chain. *arXiv preprint arXiv:1708.06733*, 2019.
- [7] Kunzhe Huang, Yiming Li, Baoyuan Wu, Zhan Qin, and Kui Ren. Backdoor defense via decoupling the training process. In *International Conference on Learning Representations*, 2022.
- [8] Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, and Xingjun Ma. Neural attention distillation: Erasing backdoor triggers from deep neural networks. In *International Conference on Learning Representations*, 2020.
- [9] Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, and Xingjun Ma. Anti-backdoor learning: Training clean models on poisoned data. *Advances in Neural Information Processing Systems*, 2021.
- [10] Yuezun Li, Yiming Li, Baoyuan Wu, Longkang Li, Ran He, and Siwei Lyu. Invisible backdoor attack with sample-specific triggers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021.
- [11] Jinlong Liu, Guoqing Jiang, Yunzhi Bai, Ting Chen, and Huayan Wang. Understanding why neural networks generalize well through gsnr of parameters. *arXiv preprint arXiv:2001.07384*, 2020.
- [12] Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. Fine-pruning: Defending against backdooring attacks on deep neural networks. In *International Symposium on Research in Attacks, Intrusions, and Defenses*, 2018.
- [13] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems*, 2017.
- [14] Tuan Anh Nguyen and Anh Tran. Input-aware dynamic backdoor attack. *Advances in Neural Information Processing Systems*, 2020.
- [15] Tuan Anh Nguyen and Anh Tuan Tran. Wanet - imperceptible warping-based backdoor attack. In *International Conference on Learning Representations*, 2021.
- [16] Fanchao Qi, Yangyi Chen, Mukai Li, Yuan Yao, Zhiyuan Liu, and Maosong Sun. ONION: A simple and effective defense against textual backdoor attacks. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 2021.
- [17] Fanchao Qi, Mukai Li, Yangyi Chen, Zhengyan Zhang, Zhiyuan Liu, Yasheng Wang, and Maosong Sun. Hidden killer: Invisible textual backdoor attacks with syntactic trigger. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing*, 2021.
- [18] Fanchao Qi, Yuan Yao, Sophia Xu, Zhiyuan Liu, and Maosong Sun. Turn the combination lock: Learnable textual backdoor attacks via word substitution. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing*, 2021.

- [19] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, 2017.
- [20] Ali Shafahi, W Ronny Huang, Mahyar Najibi, Octavian Suciuc, Christoph Studer, Tudor Dumitras, and Tom Goldstein. Poison frogs! targeted clean-label poisoning attacks on neural networks. *Advances in neural information processing systems*, 2018.
- [21] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pages 1631–1642, 2013.
- [22] Suraj Srinivas and François Fleuret. Full-gradient representation for neural network visualization. *Advances in neural information processing systems*, 32, 2019.
- [23] Mariya Toneva, Alessandro Sordani, Remi Tachet des Combes, Adam Trischler, Yoshua Bengio, and Geoffrey J. Gordon. An empirical study of example forgetting during deep neural network learning. In *7th International Conference on Learning Representations*, 2019.
- [24] Brandon Tran, Jerry Li, and Aleksander Madry. Spectral signatures in backdoor attacks. In *Advances in Neural Information Processing Systems Workshop*, 2018.
- [25] Alexander Turner, Dimitris Tsipras, and Aleksander Madry. Label-consistent backdoor attacks. *arXiv preprint arXiv:1912.02771*, 2019.
- [26] Bolun Wang, Yuanshun Yao, Shawn Shan, Huiying Li, Bimal Viswanath, Haitao Zheng, and Ben Y Zhao. Neural cleanse: Identifying and mitigating backdoor attacks in neural networks. In *2019 IEEE Symposium on Security and Privacy*, 2019.
- [27] Dongxian Wu and Yisen Wang. Adversarial neuron pruning purifies backdoored deep models. *Advances in Neural Information Processing Systems*, 2021.
- [28] Marcos Zampieri, Shervin Malmasi, Preslav Nakov, Sara Rosenthal, Noura Farra, and Ritesh Kumar. Predicting the type and target of offensive posts in social media. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2019.
- [29] Qing T Zeng, Sergey Goryachev, Scott Weiss, Margarita Sordo, Shawn N Murphy, and Ross Lazarus. Extracting principal diagnosis, co-morbidity and smoking status for asthma research: evaluation of a natural language processing system. *BMC medical informatics and decision making*, 6(1):1–9, 2006.
- [30] Yi Zeng, Won Park, Z Morley Mao, and Ruoxi Jia. Rethinking the backdoor attacks’ triggers: A frequency perspective. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021.