000 TOWARDS RELIABLE BACKDOOR ATTACKS ON 001 VISION TRANSFORMERS 002 003

Anonymous authors

Paper under double-blind review

ABSTRACT

Backdoor attacks, which make Convolution Neural Networks (CNNs) exhibit specific behaviors in the presence of a predefined trigger, bring risks to the usage of CNNs. These threats should be also considered on Vision Transformers. However, previous studies found that the existing backdoor attacks are powerful enough in ViTs to bypass common backdoor defenses, *i.e.*, these defenses either fail to reduce the attack success rate or cause a significant accuracy drop. This study investigates the existing backdoor attacks/defenses and finds that this kind of achievement is over-optimistic, caused by inappropriate adaption of defenses from CNNs to ViTs. Existing backdoor attacks can still be easily defended against with proper inheritance from CNNs. Furthermore, we propose a more reliable attack: adding a small perturbation on the trigger is enough to help existing attacks more persistent against various defenses. We hope our contributions, including the finding that existing attacks are still easy to defend with adaptations and the new backdoor attack, will promote more in-depth research into the backdoor robustness of ViTs.

023 024 025

026 027

004

010 011

012

013

014

015

016

017

018

019

021

INTRODUCTION 1

Vision Transformers (ViTs) (Dosovitskiy et al., 2021; Liu et al., Table 1: The performance 028 029 2021) have demonstrated outstanding performance in various tasks, including image classification (Yuan et al., 2021; Touvron et al., 2022), semantic segmentation (Strudel et al., 2021), and image 031 generation (Hirose et al., 2021; Bao et al., 2022), leading to their widespread popularity. However, strong performance alone is in-033 sufficient for ViT to be practically deployable. It must also exhibit 034 security and trustworthiness without posing severe security risks. One of the most notable threats to the security of ViTs is backdoor attacks (Gu et al., 2017; Chen et al., 2017), which implant unex-

of FT against Badnets attack for ResNet-18 and ViT-B on CIFAR-10 (Wu et al., 2022).

	ResNet18	ViT-B
ASR	1.48%	8.81%
ACC	89.96%	42.00%

037 pected behaviors inside models, making the victim model produce specific misclassification in the 038 presence of a predefined trigger while maintaining high performance on benign images. While previous studies mainly focus on convolution neural networks (CNNs), there is a growing need for an in-depth investigation of ViTs to help practitioners better understand the potential risks and deploy 040 them more reliably. 041

042 After a long arms race between backdoor attack and defense, for CNNs, a relatively simple defense 043 has the potential to make backdoor attacks fail, taking fine-tuning defense and Badnets attack as an 044 example in Table 1, we find that Badnets attack makes the attack success rate (ASR) on ResNet18 only have 1.48% while the benign accuracy (ACC) is 89.96%, which indicates a comprehensive failure of the attack under defense. Contrastingly, ViTs, when subjected to the same attack, display an 046 increased ASR and decreased ACC, implying the disruption of the benign utility. Given that Badnets 047 is model-agnostic, this differential outcome piqued our interest, driving us to explore the underlying 048 disparities between CNNs and ViTs. 049

Drawing inspiration from Mo et al. (2022), we discerned a crucial observation: 1) CNNs are usually 051 trained by SGD and its fine-tuning defense is also trained by SGD; 2) ViTs are typically trained by AdamW while its fine-tuning defense is trained by SGD (NOT AdamW, inheriting from earliest work 052 (Dosovitskiy et al., 2021), which first introduces optimizers to computer vision). This discrepancy in optimizers raises the possibility that the perceived vulnerability of ViTs (with defense) might be overstated, i.e., the success of attacks on ViTs with defense may be questionable. In this paper, we
 first conduct a series of experiments to comprehensively investigate the above hypothesis, which
 is further confirmed that the threat posed to ViTs with defense has been magnified. Upon minor
 modifications, ViTs with existing backdoor defense methods demonstrate clear resistance to attacks,
 mirroring the robustness of CNNs.

To this end, we are wondering whether a more powerful attack exists that can better evade current defenses. Therefore, we analyze backdoored models and further propose a simple yet effective attack. We discover that it is easy for backdoor defenses to detect and utilize the differences in channel activations due to the noticeable difference in the intermediate layers between the inputs with and without triggers. However, we can reduce this difference by adding small perturbations to the triggers before training while keeping triggers unchanged during testing, resulting in more reliable backdoor attacks. Additionally, our method has transferability across different transformer architectures and is effective for both small and large datasets.

067 068

069

071

073

075

076

077 078

079

081 082 In summary, our contributions are summarized as follows:

- We investigate the existing backdoor defenses on ViTs and find the outstanding performance of the backdoor attacks to ViTs is over-estimated due to the inappropriate adaption from CNNs to ViTs. Further, we provide a practical training recipe to improve the defense performance of existing methods and show that existing attacks can not provide reliable performances after defense.
 - We propose to add small perturbations to the triggers before training to suppress the difference in the intermediate-level representations between the inputs with and without triggers, resulting in a reliable attack. The proposed method can transfer across various architectures.
 - Our contributions, including the finding of existing attacks to current defenses and the development of a new attack, contribute to a reliable baseline for the backdoor robustness of ViTs. We hope it can be a cornerstone of future studies in the backdoor robustness of ViTs.

2 RELATED WORK

083 084 085

094

2.1 BACKDOOR ATTACK

Backdoor attacks Gu et al. (2017); Chen et al. (2017), also known as Trojan attacks, indicate the
behaviors of implanting specific malicious behavior into machine learning models, which make the
models perform well on benign data while leading to specific misclassifications on inputs containing
triggers (*i.e.*, triggered inputs). The adversary usually poisons the training data (Zeng et al., 2021) or
controls the training process (Liu et al., 2018b) to achieve this. Typically, a trigger pattern is added to
the input image as follows,

$$\boldsymbol{x}_p = (\boldsymbol{1} - \boldsymbol{m}) \odot \boldsymbol{x} + \boldsymbol{m} \odot \boldsymbol{t}, \tag{1}$$

where t is the trigger pattern and mask m indicates the pixels affected by the trigger pattern. Usually, 095 the adversary re-labels the triggered input as the predefined target class (*i.e.* in a dirty-label setting). 096 Models trained on a mixture of these poisoned data and other benign data are implanted with an unexpected correlation between the trigger pattern and the target class. To improve the stealthiness of 098 the attacks, some studies explored less noticeable trigger designs like the semi-transparent trigger (Chen et al., 2017), the elastic transformed trigger (Nguyen & Tran, 2021), and the input-aware 100 trigger (Nguyen & Tran, 2020). Besides, since incorrect annotation might expose the existence of 101 triggered data, some studies focus on poisoning without re-labeling (clean-label settings) (Turner 102 et al., 2019; Barni et al., 2019; Shafahi et al., 2018). Although most previous backdoor attacks 103 focus on CNNs, researchers have started to focus on backdoor attacks on ViT since their increasing 104 popularity. Although ViTs are reported to be more robust against adversarial attacks (Aldahdooh 105 et al., 2021; Shao et al., 2021) and common corruption (Bai et al., 2021; Bhojanapalli et al., 2021), they are still vulnerable to backdoor attacks (Lv et al., 2021; Subramanya et al., 2022). Reliable 106 attacks are needed to help practitioners properly understand the risks of backdoor attacks and deploy 107 these models reliably.

108 2.2 BACKDOOR DEFENSE

110 To mitigate the potential risks caused by backdoor attacks, numerous studies proposed various defense methods, mainly categorized into **defense during training** and **defense after training** based on the 111 stages at which they are applied. Defense during training attempts to mitigate the impact of poisoned 112 data in the training set. Some methods detect and remove poisoned data by treating them as outliers 113 (Chou et al., 2018; Udeshi et al., 2022; Gao et al., 2019), some employ semi-supervised learning 114 to bypass the incorrect correlations (Huang et al., 2022), and others utilize differential privacy to 115 ensure that a poisoned portion of training data is unable to cause severe results (Miao et al., 2022). 116 Meanwhile, the defense after training directly removes the backdoor behavior inside DNNs. This can 117 be accomplished by fine-tuning the model using a small amount of clean data (Sha et al., 2022) and it 118 can be further enhanced by first pruning the inactivated neuron (Liu et al., 2018a) or encouraging 119 the alignment of attentions (Li et al., 2021) between the student and the teacher network. Since the 120 performances of fine-tuning are easy to suffer a substantial decrease when the data is limited, another 121 popular method is selectively removing neurons related to the backdoor behaviors (Wu & Wang, 122 2021; Chai & Chen, 2022; Wang et al., 2019): Built upon the observation that the backdoor behavior can be revealed by the adversarial neuron perturbation, ANP (Wu & Wang, 2021) formulates the 123 following min-max problem with dataset D_v to expose the malicious neuron: 124

$$\min_{\mathbf{m}\in[0,1]^n} \left[\alpha \mathcal{L}_{D_v}(\mathbf{m} \odot$$

+
$$(1-\alpha) \max_{\boldsymbol{\delta}, \boldsymbol{\xi} \in [-\epsilon, \epsilon]^n} \mathcal{L}_{D_v}((\mathbf{m} + \boldsymbol{\delta}) \odot \mathbf{w}, (1+\boldsymbol{\xi})\mathbf{b})],$$

 $\mathbf{w}, \mathbf{b})$

(2)

where δ and ξ are the perturbations to the weight w and bias b of all neurons respectively. They 130 maximize the cross-entropy loss \mathcal{L}_{D_v} and **m** is the mask that adversarially preserves the clean 131 accuracy and covers up the backdoor behavior. Then the neurons corresponding to low mask values 132 are pruned to purify the backdoor model. As an improved approach based on ANP, AWM in (Chai 133 & Chen, 2022) proposes to adopt the element-wise weight masking strategies and perturb the input 134 data instead of the neurons to gain better performances on small networks. This paper primarily 135 focuses on defense after training. Because ViTs demand a large amount of data and extensive training 136 resources, it has become impractical for most practitioners to train ViTs from scratch, making defense 137 after training a more realistic scenario. Previous studies (Wu et al., 2022; Yuan et al., 2023) suggested that directly applying defenses from CNNs to ViTs fails. For example, fine-tuning decreases natural 138 accuracy from 94.58% to 42.00% against the Badnets attack and fine-pruning totally collapses in 139 (Yuan et al., 2023). At the meantime, only a few defense methods specially designed for ViT are 140 proposed (Doan et al., 2022; Subramanya et al., 2024) and their performance is lagging far behind 141 the state-of-the-art defense on CNNs: The adaptive defense proposed in (Zheng et al., 2022) only 142 decreases the ASR of TrojViT (a ViT-specific attack) to 77.13% and the patch processing method in 143 (Doan et al., 2022) fails to detect 33.2% backdoor examples on CIFAR-10. It seems that existing 144 attacks can already obtain outstanding performances on resisting defense for ViTs. However, in this 145 paper, after re-investigating various backdoor defenses with ViTs, we reveal that the achievement 146 obtained by previous attacks is not reliable. Furthermore, we provide a reliable attack, based on the 147 empirical observation of the channel activations of ViTs. It might help future research on backdoor 148 robustness with ViTs.

149

125

126

127 128 129

150 151

3 THE VULNERABILITY OF VITS (WITH DEFENSE) TO EXISTING ATTACKS

In this section, we reevaluate the perceived susceptibility of ViTs to prevailing backdoor attacks when equipped with potential defenses. We primarily consider two categories of defenses: one is fine-tuning-based, including Fine-Tuning (FT) (Sha et al., 2022), Fine-Pruning (FP) (Liu et al., 2018a), and Neural Attention Distillation (NAD) (Li et al., 2021), and the other is pruning-based, including Adversarial Neuron Pruning (ANP) (Wu & Wang, 2021) and Adversarial Weight Masking (AWM) (Chai & Chen, 2022).

158

159 3.1 BASIC SETTINGS

Here, we train a backdoored ViT-B (Dosovitskiy et al., 2021) with various attack methods. Specifically, we initialize the model with a pre-trained weight (Wightman, 2019) on the ImageNet-1k (Deng et al.,

Attack		ACC			ASR				
Tituen	No defense	SGD AdamW		No defense	SGD	AdamW			
Badnets	97.85	58.74	93.79	100.00	3.40	2.51			
Blend	97.85	94.33	93.30	100.00	13.49	4.91			
CLB	97.83	94.60	94.06	96.23	10.49	1.33			
SIG	97.50	51.56	93.51	90.57	2.23	1.40			
AvgDrop	-	22.95	4.10↓	-	89.30	94.16↑			

Table 2: The comparison between SGD and AdamW optimizer on FT. Here, AvgDrop represents the average drop of four attacks on ASR/ACC after performing FT.

173 174

2009) and then fine-tune it on CIFAR- 10^1 (Krizhevsky et al., 2009). Note that a portion of CIFAR-10175 training data is contaminated to implant the backdoor behavior, *i.e.*, some images are added with the 176 trigger pattern and are re-labeled as the target class if expected. We apply four commonly-used attack methods: 1) Badnets (Gu et al., 2019), 2) Blend (Chen et al., 2017), 3) CLB (Turner et al., 2019), and 177 4) SIG (Barni et al., 2019). Their trigger design and poisoning method in the original paper are kept. 178 To accommodate the input size of ViT, we first add triggers to CIFAR-10 images (32×32) and then 179 resize them to a larger size (224×224) . For detailed information, please refer to Appendix A. Here, 180 we use accuracy (ACC) to indicate the classification performance on benign data, and attack success 181 rate (ASR), the percentage of triggered input being classified as the target class, to indicate the attack 182 performance. Note that we will remove the inputs whose ground-truth label is the target class, and 183 thus, a successful defense should make ASR as low as 0.

184 185

3.2 VITS WITH FINE-TUNING-BASED DEFENSE

187 Fine-tuning is one of the most basic and model-agnostic defenses. However, as discussed in Section 1, 188 directly inheriting fine-tuning-based defense strategies from CNNs can potentially lead to suboptimal 189 outcomes. Note that SGD is the commonly used optimizer for both training and fine-tuning for 190 CNNs, while for ViTs, the first work (Dosovitskiy et al., 2021) introducing Transformers to computer 191 vision, adopts AdamW for pre-training and SGD for fine-tuning. Notably, prior work (Wu et al., 192 2022) on backdoor defense naturally inherit this strategy and observes notably diminished accuracy 193 across multiple backdoor attacks. This discrepancy in optimizers motivates us to study the potential influence of optimizers on backdoor defense. The initial learning rates for SGD and AdamW are set 194 to 0.02 and 3e-4, respectively. For the other parameters in AdamW, we use the common settings of the 195 original ViTs (refer to Appendix B for details). Table 2 illustrates the experimental fine-tuning (FT) 196 results against various backdoor attacks. For the results on FP and NAD, please refer to Appendix C. 197 We find that SGD exhibited significant instability on ViTs. Even for the same model, when defending 198 against Blend and CLB, it achieves more than 90% of ACC. However, for BadNet and SIG, ACC 199 decreases to less than 60%. In contrast, AdamW consistently achieves high ACC and low ASR using 200 the same hyper-parameter configuration. Therefore, simply using SGD for backdoor defense on ViTs 201 will yield highly unstable performance. We recommend employing AdamW for defense purposes.

202 203

204

3.3 VITS WITH PRUNING-BASED DEFENSE

205 Pruning is also a typical defense approach, which attempts to remove backdoor-related neu-206 rons/channels and is severely impacted by the architectures. In previous studies, pruning-based 207 methods have achieved excellent robustness against backdoor attacks with CNNs (Wu & Wang, 2021; 208 Chai & Chen, 2022). However, when we directly apply these methods to ViTs, we find that they 209 are unable to effectively defend as shown in Table 3. Specifically, ANP fails to reduce ASR and 210 cannot remove the backdoor-related neurons. Besides, although AWM reduces ASR, it also severely decreases ACC, making the model unusable. To explore the potential reason, we look deeply at the 211 implementation of ANP and find that ANP actually prunes channels inside norm layers rather than 212 neurons inside convolutional layers. This is because, in CNNs, each neuron is typically surrounded 213

¹⁷¹ 172

¹Ony 95% of the original training data on CIFAR-10 are used to train the backdoored model, and the remaining data are kept for defense.

Metric	Setting	Before	ANP	ANP (Adapted)	AWM	AWM (Adapted)
	Badnets	97.85	97.85	94.26	85.98	95.02
ACC	Blend	97.85	97.85	92.70	83.29	95.08
ACC	CLB	97.83	97.83	95.71	85.67	95.60
	SIG	97.50	97.50	92.60	87.22	94.58
	Badnets	100.00	100.00	1.34	1.24	0.71
ACD	Blend	100.00	100.00	23.7	2.03	1.70
ASK	CLB	96.23	96.23	12.71	3.48	1.52
	SIG	90.57	90.57	1.48	1.16	3.87

Table 3: The Performance of pruning-based defense with or without ViTs adaptation.



Figure 1: The average activations for different channels before (a) and after the backdoor defense (b)-(c). The activations are sorted in descending order of the activations on natural samples.

by at least one norm layer². However, in ViT, many norm layers are removed, and norm-layer-based pruning only influences part of neurons and limits the defense performance. Meanwhile, AWM utilizes element-wise masks for optimization, whose number of parameters is the same as the total number of parameters of ViT. Since ViTs are typically larger, AWM encounters the severe overfitting issue, leading to low accuracy. Therefore, to make pruning methods applicable to ViTs, selecting appropriate granularity and pruning locations is necessary. Here, we recommend directly pruning all channels of linear projection inside both attention and MLP layers, which provides better coverage than ANP and requires fewer parameters compared to AWM. This modification decreases ASR notably and keeps ACC high.

4 PROPOSED BACKDOOR ATTACKS

Following the above analysis, existing defense methods (ViTs adapted) successfully defend against
backdoor attacks in ViTs, just as they do in CNNs. Here, we want to explore whether there exist new
backdoor attacks to beat the newly adapted defense on ViTs.

To obtain a better insight into why defense methods can detect and remove backdoor behaviors, we investigate the per-channel activations before the MLP head in ViT. We illustrate the average activations of all channels for a backdoored ViT-B on triggered and benign inputs from the CIFAR-10 test set, respectively. For clarity, we reorganize the channels based on their average activations, arranging them from largest to smallest with respect to average activations on benign data. In Figure 1, we find a significant activation difference between benign and triggered inputs, which is easy to capture. Further, we compare the average activation of all channels for models purified by FT and AWM, and find that benign and triggered inputs have similar average activation after defense. This suggests that the naive trigger design (usually predefined universal patterns) for current backdoor attacks results in a significant difference between benign and triggered data, revealing attack information to possible defenders. Next, we will study whether we could improve the trigger design

²Specifically, for Preact-ResNet, the norm layer is always located before the neuron; for ResNet, it is located after the neuron



Figure 2: The illustration of our proposed attack. We illustrate our attack by taking ViT-B as an example. *left:* Using the existing poisoned dataset, we also train the BD and TC simultaneously during backdoor training (**Step 1**). *right:* When the training is over, we perform adversarial attacks on the BD and TC modules to generate adversarial perturbation (**Step 2**). In each step during crafting adversarial perturbations, we manually mask some patches of perturbation to better poisoned ViTs (**Step 3**).

to escape defenses. The general process of our attack is summarized in Figure 2 and we term it as the Channel Activation attack in ViT (CAT).

Adversarial Loss. Based on our observation, a good trigger design is expected to avoid noticeable 298 channel activation differences between benign and triggered inputs. Therefore, we require additional 299 backdoor discriminators (BD) to clarify whether the training input has the predefined trigger during 300 the training. Specifically, we denote the feature extractor of the backdoored model as $q(\cdot)^3$, and 301 the backdoor discriminator $d_i(q(x))$ uses the intermediate feature of the *i*-th layer to discriminate 302 whether the input x has the trigger pattern. During backdoor training, we also train these backdoor 303 discriminators of the last n layers, *i.e.*, $d_i(g(x)), i = L - n + 1, \dots, L$. After training, we could 304 use these backdoor discriminators to generate adversarial perturbations on the trigger pattern to 305 minimize the activation difference between benign and triggered inputs. Meanwhile, naive difference minimization might make the model classify triggered inputs as a non-target label, leading to 306 the failure of backdoor attacks. To address this issue, we introduce additional target classifiers 307 $f_i(q(x))$ (TC), which uses the intermediate feature of the *i*-th layer to make classification between 308 benign samples, *i.e.*, classifying the benign input as the ground-truth label. Similar to the backdoor 309 discriminator, we also train these clean classifiers of the last n layers, *i.e.*, $f_i(g(x)), i = L - L$ 310 $n + 1, \dots, L$ during training. In conclusion, we craft adversarial perturbation via maximizing the 311 following loss, L 312

$$\mathcal{L}(\boldsymbol{\delta}) = \sum_{i=L-n+1}^{-} (1-\gamma) \cdot \ell \big(d_i (\boldsymbol{g}(\boldsymbol{x} + \boldsymbol{m} \odot \boldsymbol{\delta})), y_{\text{bd}} \big) - \gamma \cdot \ell \big(f_i (\boldsymbol{g}(\boldsymbol{x} + \boldsymbol{m} \odot \boldsymbol{\delta})), y_{\text{tc}} \big),$$
(3)

where y_{bd} is the label for the backdoor discriminator, *i.e.*, 1 for triggered data and 0 for benign data. y_{tc} is the label for the target classifier as the adversary expects, *i.e.*, the ground-truth label for benign input, and the target label for triggered input. Here γ is a trade-off coefficient to balance the effect between TC and BD.

Generation Steps. Since the nonlinearity of ViTs, it is mathematically infeasible to obtain the exact solution for Equation 3. However, we can use the projected gradient descent (PGD) (Madry et al.,

323

288

289

290

291

292

293

295

296

297

L

³In our method, the extractor will return intermediate features from all layers.

2018) from the normal adversarial attacks to craft the perturbations on the trigger pattern as follows: 325

$$\boldsymbol{\delta} \leftarrow \boldsymbol{m} \odot \Pi_{\epsilon} \big(\boldsymbol{\delta} + \alpha \cdot \frac{\nabla_{\boldsymbol{\delta}} \mathcal{L}(\boldsymbol{\delta})}{\|\nabla_{\boldsymbol{\delta}} \mathcal{L}(\boldsymbol{\delta})\|_2} \big), \tag{4}$$

where *m* is the mask for triggers, \odot is the Hadamard product, and $\Pi_{\epsilon}(\cdot)$ is the projection function,

$$\Pi_{\epsilon}(\boldsymbol{\delta}) = \frac{\epsilon}{\|\boldsymbol{\delta}\|_2} \boldsymbol{\delta}.$$
(5)

Random Masking of Perturbation. In practical situations, the adversary has no access to model architecture and its parameters. Usually, the adversary expects to craft these perturbations from models with known parameters and structure (source model) to attack these unknown models (target model). The generated perturbations in this situation are expected to be effective across various architectures. Unfortunately, different ViTs could have various patch sizes for splitting, leading to differences in the scale of sensitive features. This might cause low transferability across architectures. Therefore, we propose a method termed Random Masking of Perturbation (RMP). In each step during crafting adversarial perturbations, we first split perturbation with k patches and randomly drop a predefined percentage of perturbation patches. This can create features of varying scales manually and make the perturbations effective for kinds of ViTs with different patch-splitting approaches.

341 342 343

344

346

324

326 327 328

330 331 332

333

334

335

336

337

338

339

340

- 5 **EXPERIMENTS**
- 345 5.1 MAIN RESULTS

Settings: We evaluate the performances of our methods in two scenarios. 1) White-box: the target 347 model and source models have the same architectures and backdoor training from the same pre-trained 348 model. 2) Black-box: the architectures of the target model and the source model are different. We 349 choose ViT-B as the source model and five ViT variants, including ViT-B, DeiT-S (Touvron et al., 350 2021a), Swin-B (Liu et al., 2021), Cait-S (Touvron et al., 2021b) and XciT-S (Ali et al., 2021) as 351 our target models. In our experiments, we choose the last two layers (*i.e.*, n = 2) to add BD and 352 TC modules, which are composed of one Linear layer. For the perturbation generation step, the 353 adversarial attack is l_2 bounded PGD-10 with budget 16/255, step size 4/255, and the trade-off 354 parameter γ is set to 0.6. For random masking of perturbation, we split the perturbation into multiple 355 small pieces, each of which has the shape of 2×2 . The percentage of dropped patches is set to 0.1 356 and 0.05 for the whole-image patch and trigger-based path, respectively. For other hyperparameters, we mainly keep in line with Section 3 and summarize them in Appendix A and B. All experiments 357 are performed on CIFAR-10. The ASR of our CAT against five defenses are summarized in Table 358 4. For ACC, please refer to Appendix D. In addition to the five defenses mentioned in previous 359 sections, in Appendix E, we also demonstrate that CAT can even help prevailing attacks bypass the 360 detected-based defenses. 361

Results: First, when no defenses are performed, CAT will obtain a comparable ASR compared to 362 the vanilla settings. In most cases, it even can gain better performance. For example, our method increases the ASR of SIG attack from 90.57% to 91.19% on ViT-B. Second, for the post-defense 364 situation, CAT can achieve higher ASR in a novel margin. For example, under the white-box setting, it increases the ASR from 2.51% to 66.72% against the badnets attack for FT. In the black-box 366 settings, the ASR of SIG attacks increases from 3.30% to 13.81% on DeiT-S for the AWM defenses. 367 As for ACC, the results in Appendix D show that CAT will obtain comparable ACC compared to the 368 vanilla attack. It indicates that our method will only enhance the ASR without compromising the 369 classification of the benign images.

370 371

372

5.2 PERFORMANCE ON IMAGENET WITH COMPARISONS WITH VIT-SPECIFIC METHODS

373 Attribute to the highly flexible multi-head self-attention mechanism, ViTs can outperform CNNs 374 when millions of data are provided. Thus in this section, we not only evaluate the performance of 375 our attack on ImageNet (Deng et al., 2009) but also compare it with existing ViT-specific attacks to illustrate its superiority. Here we only report the results after combining badnets and blend attacks 376 because the clean-label attacks will fail for only at-most poisoning 0.1% of training data. More details 377 of our experimental configurations are summarized in Appendix F. In addition to the model-agnostic

Definite Vanilla CAT Vanilla Vanilla CAT	Defense	Attack	Vi	Г-В	Dei	T-S	Swi	n-B	Cai	T-S	Xci	T-S
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	Derense	. maion	Vanilla	CAT								
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $		BadNets	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	No defense	Blend	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	No defense	CLB	96.23	94.57	95.28	94.04	84.86	90.23	85.71	92.21	100.00	100.00
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $		SIG	90.57	91.19	84.77	88.28	94.99	97.77	80.93	82.26	94.21	96.17
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $		BadNets	2.51	66.72	8.21	56.82	8.17	36.07	16.57	64.34	8.81	51.80
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	FT	Blend	4.91	38.53	1.84	18.94	1.84	15.36	10.11	49.12	11.58	85.93
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	Defense No defense FT FP NAD ANP AWM	CLB	1.33	12.32	6.39	18.39	0.90	8.22	0.48	10.93	9.88	63.39
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $		SIG	1.40	10.99	0.88	4.20	0.18	15.02	14.54	19.82	3.09	25.18
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $		BadNets	0.91	27.90	33.98	45.09	11.49	19.52	15.80	19.96	6.37	14.39
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	ED	Blend	0.73	12.49	3.82	14.73	2.48	22.67	43.09	90.27	23.82	29.50
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	1.1	CLB	1.70	26.88	3.87	16.52	2.54	5.56	1.59	6.12	13.99	20.49
BadNets 1.57 86.50 10.20 43.73 4.26 47.09 3.83 32.52 15.20 32.63 NAD Blend 8.94 61.93 1.50 23.30 5.70 59.32 11.89 73.71 18.57 55.90 SIG 7.27 13.30 5.88 19.66 1.32 11.61 3.84 29.66 20.10 36.04 SIG 3.60 9.07 3.99 21.02 1.27 24.62 4.39 18.17 6.16 23.01 ANP Blend 23.70 92.23 36.67 79.91 34.37 99.62 59.83 100.00 0.00 99.99 ANP Blend 12.71 14.01 13.18 25.19 2.78 10.62 2.64 23.51 44.53 92.43 SIG 1.48 67.57 20.8 79.91 21.78 60.26 41.88 67.63 21.72 64.52 AWM		SIG	0.81	9.68	2.26	14.79	3.81	5.49	8.96	19.23	13.22	16.43
NAD Blend 8.94 61.93 1.50 23.30 5.70 59.32 11.89 73.71 18.57 55.90 CLB 7.27 13.30 5.88 19.66 1.32 11.61 3.84 29.66 20.10 36.04 SIG 3.60 9.07 3.99 21.02 1.27 24.62 4.39 18.17 6.16 23.01 ANP Blend 23.70 92.23 36.67 79.91 34.37 99.62 59.83 100.00 0.00 99.99 CLB 12.71 14.01 13.18 25.19 2.78 10.62 2.64 23.51 44.53 92.43 SIG 1.48 67.57 20.8 79.91 21.78 10.62 2.64 23.51 44.53 92.43 SIG 1.48 67.57 20.8 79.91 21.76 0.00 10.57 2.31 16.11 AWM Blend 1.70		BadNets	1.57	86.50	10.20	43.73	4.26	47.09	3.83	32.52	15.20	32.63
NAD CLB 7.27 13.30 5.88 19.66 1.32 11.61 3.84 29.66 20.10 36.04 SIG 3.60 9.07 3.99 21.02 1.27 24.62 4.39 18.17 6.16 23.01 ANP BadNets 1.34 51.09 6.03 58.17 2.63 19.47 31.24 83.34 6.82 81.57 ANP Blend 23.70 92.23 36.67 79.91 34.37 99.62 2.64 23.51 44.53 92.43 SIG 1.48 67.57 20.8 79.91 21.78 60.26 41.88 67.63 21.72 64.52 AWM Blend 1.70 26.22 1.27 5.12 0.32 27.62 36.00 57.72 88.43 94.56 AWM Blend 1.70 26.22 1.27 5.12 0.32 27.62 36.00 57.72 88.43 94.56 SIG 3.87	No defense FT FP NAD ANP AWM	Blend	8.94	61.93	1.50	23.30	5.70	59.32	11.89	73.71	18.57	55.90
SIG 3.60 9.07 3.99 21.02 1.27 24.62 4.39 18.17 6.16 23.01 ANP BadNets 1.34 51.09 6.03 58.17 2.63 19.47 31.24 83.34 6.82 81.57 ANP Blend 23.70 92.23 36.67 79.91 34.37 99.62 59.83 100.00 0.00 99.99 CLB 12.71 14.01 13.18 25.19 2.78 10.62 2.64 23.51 44.53 92.43 SIG 1.48 67.57 20.8 79.91 21.78 60.26 41.88 67.63 21.72 64.52 AWM Blend 1.70 26.72 1.271 5.12 0.32 27.62 36.00 57.72 88.43 94.56 AWM Blend 1.70 26.22 1.27 5.12 0.32 27.62 36.00 57.72 88.43 94.55 SIG 3.87 38.59		CLB	7.27	13.30	5.88	19.66	1.32	11.61	3.84	29.66	20.10	36.04
BadNets 1.34 51.09 6.03 58.17 2.63 19.47 31.24 83.34 6.82 81.57 ANP Blend 23.70 92.23 36.67 79.91 34.37 99.62 59.83 100.00 0.00 99.99 CLB 12.71 14.01 13.18 25.19 2.78 10.62 2.64 23.51 44.53 92.43 SIG 1.48 67.57 20.8 79.91 21.78 60.26 41.88 67.63 21.72 64.52 AWM BadNets 0.71 6.78 2.71 6.64 4.79 12.76 0.90 10.57 2.31 16.11 AWM Blend 1.70 26.22 1.27 5.12 0.32 27.62 36.00 57.72 88.43 94.56 CLB 1.52 4.40 2.19 5.42 3.16 6.74 0.91 2.66 26.84 40.71 SIG 3.87 38.59 3.30 <		SIG	3.60	9.07	3.99	21.02	1.27	24.62	4.39	18.17	6.16	23.01
ANP Blend 23.70 92.23 36.67 79.91 34.37 99.62 59.83 100.00 0.00 99.99 CLB 12.71 14.01 13.18 25.19 2.78 10.62 2.64 23.51 44.53 92.43 SIG 1.48 67.57 20.8 79.91 21.78 60.26 41.88 67.63 21.72 64.52 AWM Blend 1.70 26.22 1.27 5.12 0.32 27.62 36.00 57.72 28.84 94.56 CLB 1.52 4.40 2.19 5.12 0.32 27.62 36.00 57.72 88.43 94.56 GLB 1.52 4.40 2.19 5.12 0.32 27.62 36.00 57.72 88.43 94.56 GLB 1.52 4.40 2.19 5.42 3.16 6.74 0.91 2.66 26.84 40.71 SIG 3.87 38.59 3.30 13.81 29		BadNets	1.34	51.09	6.03	58.17	2.63	19.47	31.24	83.34	6.82	81.57
AIM CLB 12.71 14.01 13.18 25.19 2.78 10.62 2.64 23.51 44.53 92.43 SIG 1.48 67.57 20.8 79.91 21.78 60.26 41.88 67.63 21.72 64.52 Baloets 0.71 6.78 2.71 6.64 4.79 12.76 0.90 10.57 2.31 16.17 AWM Blend 1.70 26.22 1.27 5.12 0.32 27.62 36.00 57.72 88.43 94.56 CLB 1.52 4.40 2.19 5.42 3.16 6.74 0.91 2.66 26.84 40.71 SIG 3.87 38.59 3.30 13.81 29.83 59.82 16.79 23.22 35.99 96.05	AND	Blend	23.70	92.23	36.67	79.91	34.37	99.62	59.83	100.00	0.00	99.99
SIG 1.48 67.57 20.8 79.91 21.78 60.26 41.88 67.63 21.72 64.52 BadNets 0.71 6.78 2.71 6.64 4.79 12.76 0.90 10.57 2.31 16.11 AWM Blend 1.70 26.22 1.27 5.12 0.32 27.62 36.00 57.72 88.43 94.56 CLB 1.52 4.40 2.19 5.42 3.16 6.74 0.91 2.66 2.6.84 40.71 SIG 3.87 38.59 3.30 13.81 29.83 59.82 16.79 23.22 35.99 96.05	AINT	CLB	12.71	14.01	13.18	25.19	2.78	10.62	2.64	23.51	44.53	92.43
BadNets 0.71 6.78 2.71 6.64 4.79 12.76 0.90 10.57 2.31 16.11 AWM Blend 1.70 26.22 1.27 5.12 0.32 27.62 36.00 57.72 88.43 94.56 CLB 1.52 4.40 2.19 5.42 3.16 6.74 0.91 2.66 26.84 40.71 SIG 3.87 38.59 3.30 13.81 29.83 59.82 16.79 23.22 35.99 96.05		SIG	1.48	67.57	20.8	79.91	21.78	60.26	41.88	67.63	21.72	64.52
AWM Blend 1.70 26.22 1.27 5.12 0.32 27.62 36.00 57.72 88.43 94.56 CLB 1.52 4.40 2.19 5.42 3.16 6.74 0.91 2.66 26.84 40.71 SIG 3.87 38.59 3.30 13.81 29.83 59.82 16.79 23.22 35.99 96.05		BadNets	0.71	6.78	2.71	6.64	4.79	12.76	0.90	10.57	2.31	16.11
CLB 1.52 4.40 2.19 5.42 3.16 6.74 0.91 2.66 26.84 40.71 SIG 3.87 38.59 3.30 13.81 29.83 59.82 16.79 23.22 35.99 96.05	No defense FT FP NAD ANP AWM	Blend	1.70	26.22	1.27	5.12	0.32	27.62	36.00	57.72	88.43	94.56
SIG 3.87 38.59 3.30 13.81 29.83 59.82 16.79 23.22 35.99 96.05	AW M	CLB	1.52	4.40	2.19	5.42	3.16	6.74	0.91	2.66	26.84	40.71
		SIG	3.87	38.59	3.30	13.81	29.83	59.82	16.79	23.22	35.99	96.05

Table 4: ASR (%) of our proposed attack with different ViT variants on the CIFAR-10 dataset. The best results are in **bold**.

Table 5: ASR (%) of our attack on ImageNet dataset. The higher ASR is in **bold**.

Attack	Before	FT	FP	NAD	ANP	AWM	AB
TrojViT	91.08	0.14	0.11	0.16	0.46	0.18	-
DBIA	99.58	0.09	0.07	0.10	0.10	0.05	-
Badnets Badnets+CAT	100.00 100.00	27.75 51.35	3.67 14.17	26.82 28.75	18.30 44.36	24.32 81.98	3.84 12.76
Blend Blend+CAT	100.00 100.00	18.44 27.83	1.01 3.17	6.71 13.44	19.79 48.49	39.63 71.29	100.00 100.00

attacks mentioned in the previous sections, we also include two ViT-specific attacks: the Trojan 409 Insertion attack in ViT (TrojViT) (Zheng et al., 2022) and the Data-free Backdoor Injection Attack 410 (DBIA) (Lv et al., 2021) for comparison. We also evaluate our methods on the current ViT-specific 411 defense, including Attention Blocking (AB) (Subramanya et al., 2024). The hyperparameter settings 412 of ViT-specific attacks or defenses are the same as those in the original paper (Please refer to Appendix 413 F for details). Considering both the white-box and black-box attacks, ViT-B is chosen as the source 414 model and our target models include ViT-B and Swin-B. The ASR and ACC of the white-box setting 415 are summarized in Table 5 and Appendix G respectively. For the performances of CAT on the 416 black-box settings, please refer to Appendix H for the results.

417 First, similar to the results on CIFAR-10, the results reveal that CAT can help existing attacks better 418 bypass the adapted defenses. For example, our approach boosts the ASR of Badnets from 24.32% to 419 81.98% after applying AWM. In addition, compared to the existing ViT-specific backdoor attacks, 420 our method also shows its superior performance: Both TrojViT and DBIA only obtain less than 1%421 ASR after performing the adapted defense which is quite lower than those of CAT. In addition, for 422 ViT-specific defense, our method also obtains better performance: the gains on ASR are observed 423 after combining Badnets with CAT. We conjecture this is because our attack reduces the anomalous behavior of backdoor samples on ViTs by introducing benign features. This increases the difficulty 424 of detecting them from the poison dataset. AB totally fails to defend Blend or CAT+Blend because it 425 only masks a patch of images which will be less effective when encountering the whole-image attack, 426 i.e. Blend. 427

428

430

380 381 382

396 397

429 5.3 ABLATION STUDY

For our proposed CAT, there are two key components: one is to perform adversarial attacks on triggers (PA), and the other is to randomly mask patches of perturbation (RMP). To evaluate the contribution



Table 6: The ASR for different combinations of our technique. The better result is in **bold**.



of each component, we test the performances under three combinations: 1) the vanilla backdoor attacks, 2) backdoor attacks with PA, 3) backdoor attacks with both PA and RMP. Considering both white-box and black-box settings, we select ViT-B and Swin-B as the target models. We select FP and AWM to evaluate the performances of backdoor attacks since they show the most promising performances in Table 4. Other configurations are the same as those in section 5.1. We summarize the ASR for all combinations in Table 6. It reveals that PA can improve the ASR for both ViT-B and Swin-B but applying PA and RMP together can gain higher ASR. For example, under the white-box setting, the gain of PA for FP against badnets attack is 13.63%, performing PA and RMP both can further improve the ASR by 26.99%. Similar results are also observed for the black-box settings.

5.4 Hyperparameter Analysis

⁴⁶⁶ In this section, we test the effect of hyperparameters on our proposed methods. Taking Badnets ⁴⁶⁷ attacks as an example, we report the ASR after performing fine-tuning (FT) for ViT-B and Swin-B.

Attack budget: Recalling that in Section 4, we craft the adversarial samples to reduce the differences in features between the backdoor and benign data. The previous works reveal that the strength of the attacks plays a vital significance in the adversarial region. Therefore, we first investigate the effect of the attack strength ϵ on the performance of our method. As shown in Figure 3 (a), the ASR of our method increases when we increase the budget. This is because more and more features on the triggers that mismatches the benign data are removed. However, when the attack is too strong $(\epsilon > 16/255)$, the performance of our method will decrease because it makes it too hard for the network to learn backdoor information from the data.

Trade-off coefficient: γ is another important hyperparameter for our method. As shown in Figure 3 (b), the results illustrate that the adversarial information from both additional modules: the backdoor discriminator and the target classifier can improve the ASR ($\gamma = 0$ or 1.0). However, mixing the information from both can gain better performance. When $\gamma = 0.6$, our method achieves the best performance by simultaneously enhancing the information of the target class while eliminating the irrelevant features on the triggers.

- 5.5 A CLOSER LOOK AT CAT
- **Time costs:** As proposed in Section 4, CAT only increases the cost before training, thus its additional overhead is proportional to the number of poisoned samples and the attack steps for crafting adversar-



ial perturbations. We also further perform experiments on a single RTX3090 to demonstrate that CAT only brings negligible cost to practical use on ViT-B. As shown in Table 7, the preprocessing process of CAT can be completed in a few minutes. Even on large datasets such as ImageNet, the overall costs are less than 4 minutes. The results demonstrate that it is affordable to perform CAT for most attackers, which could potentially increase its threat to current AI systems.

502 Channel activations: We visualize the activation difference with or without performing CAT against 503 the badnets attack in Figure 4 and sort them in descending order. Here the activation difference refers 504 to the absolute value of the difference in activation between clean and backdoor samples in different 505 channels. For the results of Blend, CLB and SIG attacks, please refer to Appendix I for more details. 506 Compared to vanilla attacks, the results demonstrate that CAT can largely reduce the differences between the activations of the backdoor and clean samples. Therefore it effectively increases the 508 stealthiness of the combined attacks.

509 510

507

496 497

498

499

500

501

CONCLUSION 6

511 512

513 In this paper, we conduct a comprehensive evaluation of backdoor methods on ViTs and show that 514 the illustration of success achieved by current attacks to ViTs is due to inappropriate adaption of 515 defense from CNNs to ViTs. We further provide some training recipes to correctly evaluate the attack, including using AdamW rather than SGD and selecting appropriate granularity for pruning. 516 Our results demonstrate that existing attacks can not provide reliable performance after defense. 517 Therefore, we investigate why the defense method easily removes backdoor behavior and find a 518 huge difference in channel activation in intermediate layers with commonly used predefined triggers. 519 Inspired by this, we propose a more reliable attack by adding special adversarial perturbations into 520 the trigger pattern to avoid noticeable channel activation differences between benign and triggered 521 input. We hope our method, including the proposed recipes in ViTs and the new attack method, could 522 be a cornerstone of future studies on the backdoor robustness of ViTs.

523 524

525 ETHICS STATEMENT

526 527

528

529

530

531

532

The popular use of ViTs in multiple vision tasks makes us notice their security concerns and one of those is backdoor attacks. In this paper, we not only make adaptions for the existing backdoor defenses but also propose a new backdoor attack based on the differences in dimensional activations. Our contributions may help the community reliably evaluate the backdoor robustness of ViTs and the safer application of ViTs in real-world scenarios. In the meantime, the negative impact can not be simply ignored: our proposed attack could be exploited by malicious attackers to build more powerful backdoor attacks for ViTs.

533 534

REPRODUCIBILITY STATEMENT

536 537

We perform all experiments on publicly available datasets. To ensure the reproducibility of the paper, 538 we summarize the basic settings in Section 3.1 and Section 5.1. For others, we list them in Appendix A and B respectively. We will release the open-source code upon acceptance.

540	References
541	

542	Samira Abnar and Willem Zuidema	. Ouantifying attention	flow in	transformers.	In <i>arXiv</i>	. 2020.
. 1997						

- Ahmed Aldahdooh, Wassim Hamidouche, and Olivier Deforges. Reveal of vision transformers robustness against adversarial attacks. In *arXiv*, 2021.
- Alaaeldin Ali, Hugo Touvron, Mathilde Caron, Piotr Bojanowski, Matthijs Douze, Armand Joulin,
 Ivan Laptev, Natalia Neverova, Gabriel Synnaeve, Jakob Verbeek, et al. Xcit: Cross-covariance
 image transformers. In *NeurIPS*, 2021.
- Yutong Bai, Jieru Mei, Alan L Yuille, and Cihang Xie. Are transformers more robust than cnns? In *NeurIPS*, 2021.
- Fan Bao, Chongxuan Li, Yue Cao, and Jun Zhu. All are worth words: a vit backbone for score-based diffusion models. In *arXiv*, 2022.
- Mauro Barni, Kassem Kallas, and Benedetta Tondi. A new backdoor attack in cnns by training set corruption without label poisoning. In *ICIP*, 2019.
- Srinadh Bhojanapalli, Ayan Chakrabarti, Daniel Glasner, Daliang Li, Thomas Unterthiner, and
 Andreas Veit. Understanding robustness of transformers for image classification. In *ICCV*, 2021.
- Shuwen Chai and Jinghui Chen. One-shot neural backdoor erasing via adversarial weight masking. In *NeurIPS*, 2022.
- Xinyun Chen, Chang Liu, Bo Li, Kimberly Lu, and Dawn Song. Targeted backdoor attacks on deep learning systems using data poisoning. In *arXiv*, 2017.
- Edward Chou, Florian Tramèr, Giancarlo Pellegrino, and Dan Boneh. Sentinet: Detecting physical attacks against deep learning systems. In *arXiv*, 2018.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
 hierarchical image database. In *CVPR*, 2009.
- Khoa D Doan, Yingjie Lao, Peng Yang, and Ping Li. Defending backdoor attacks on vision transformer via patch processing. In *arXiv*, 2022.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image
 is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021.
- Yansong Gao, Change Xu, Derui Wang, Shiping Chen, Damith C Ranasinghe, and Surya Nepal.
 Strip: A defence against trojan attacks on deep neural networks. In *ACSA*, 2019.
- Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Identifying vulnerabilities in the
 machine learning model supply chain. In *arXiv*, 2017.
- Tianyu Gu, Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Evaluating backdooring attacks on deep neural networks. *IEEE Access*, 2019.
- Shota Hirose, Naoki Wada, Jiro Katto, and Heming Sun. Vit-gan: Using vision transformer as discriminator with adaptive data augmentation. In *ICCCI*, 2021.
- Kunzhe Huang, Yiming Li, Baoyuan Wu, Zhan Qin, and Kui Ren. Backdoor defense via decoupling the training process. In *arXiv*, 2022.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. *Tech Report*, 2009.
- Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, and Xingjun Ma. Neural attention
 distillation: Erasing backdoor triggers from deep neural networks. In *arXiv*, 2021.
- 593 Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. Fine-pruning: Defending against backdooring attacks on deep neural networks. In *RAID*, 2018a.

594 595 596	Yingqi Liu, Shiqing Ma, Yousra Aafer, Wen-Chuan Lee, Juan Zhai, Weihang Wang, and Xiangyu Zhang. Trojaning attack on neural networks. In NDSS, 2018b.
598 597 598	Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In <i>ICCV</i> , 2021.
599 600	Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In ICLR, 2018.
601 602	Peizhuo Lv, Hualong Ma, Jiachen Zhou, Ruigang Liang, Kai Chen, Shengzhi Zhang, and Yunfei Yang. Dbia: Data-free backdoor injection attack against transformer networks. In <i>arXiv</i> , 2021.
603 604	Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In <i>ICLR</i> , 2018.
605 606 607	Lu Miao, Wei Yang, Rong Hu, Lu Li, and Liusheng Huang. Against backdoor attacks in federated learning with differential privacy. In <i>ICASSP</i> , 2022.
608 609	Yichuan Mo, Dongxian Wu, Yifei Wang, Yiwen Guo, and Yisen Wang. When adversarial training meets vision transformers: Recipes from training to architecture. In <i>NeurIPS</i> , 2022.
610 611	Anh Nguyen and Anh Tran. Wanet-imperceptible warping-based backdoor attack. In arXiv, 2021.
612	Tuan Anh Nguyen and Anh Tran. Input-aware dynamic backdoor attack. In NeurIPS, 2020.
614 615	Zeyang Sha, Xinlei He, Pascal Berrang, Mathias Humbert, and Yang Zhang. Fine-tuning is all you need to mitigate backdoor attacks. In <i>arXiv</i> , 2022.
616 617 618	Ali Shafahi, W Ronny Huang, Mahyar Najibi, Octavian Suciu, Christoph Studer, Tudor Dumitras, and Tom Goldstein. Poison frogs! targeted clean-label poisoning attacks on neural networks. In <i>NeurIPS</i> , 2018.
620 621	Rulin Shao, Zhouxing Shi, Jinfeng Yi, Pin-Yu Chen, and Cho-Jui Hsieh. On the adversarial robustness of visual transformers. In <i>arXiv</i> , 2021.
622 623 624	Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Segmenter: Transformer for semantic segmentation. In <i>ICCV</i> , 2021.
625 626	Akshayvarun Subramanya, Aniruddha Saha, Soroush Abbasi Koohpayegani, Ajinkya Tejankar, and Hamed Pirsiavash. Backdoor attacks on vision transformers. In <i>arXiv</i> , 2022.
627 628 629	Akshayvarun Subramanya, Soroush Abbasi Koohpayegani, Aniruddha Saha, Ajinkya Tejankar, and Hamed Pirsiavash. A closer look at robustness of vision transformers to backdoor attacks. In <i>WACV</i> , 2024.
630 631 632	Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In <i>ICML</i> , 2021a.
633 634 625	Hugo Touvron, Matthieu Cord, Alexandre Sablayrolles, Gabriel Synnaeve, and Hervé Jégou. Going deeper with image transformers. In <i>ICCV</i> , 2021b.
636	Hugo Touvron, Matthieu Cord, and Hervé Jégou. Deit iii: Revenge of the vit. In ECCV, 2022.
637 638 620	Alexander Turner, Dimitris Tsipras, and Aleksander Madry. Label-consistent backdoor attacks. In <i>arXiv</i> , 2019.
640 641 642	Sakshi Udeshi, Shanshan Peng, Gerald Woo, Lionell Loh, Louth Rawshan, and Sudipta Chattopad- hyay. Model agnostic defence against backdoor attacks in machine learning. <i>IEEE Transactions</i> <i>on Reliability</i> , 2022.
643 644 645 646	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 <i>S&P</i> , 2019.
647	Ross Wightman. Pytorch image models. https://github.com/rwightman/ pytorch-image-models, 2019.

- Baoyuan Wu, Hongrui Chen, Mingda Zhang, Zihao Zhu, Shaokui Wei, Danni Yuan, and Chao Shen.
 Backdoorbench: A comprehensive benchmark of backdoor learning. In *NeurIPS*, 2022.
- Dongxian Wu and Yisen Wang. Adversarial neuron pruning purifies backdoored deep models. In
 NeurIPS, 2021.
- Dongxian Wu, Shu-Tao Xia, and Yisen Wang. Adversarial weight perturbation helps robust general ization. In *NeurIPS*, 2020.
- Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Zi-Hang Jiang, Francis EH Tay, Jiashi
 Feng, and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on
 imagenet. In *ICCV*, 2021.
 - Zenghui Yuan, Pan Zhou, Kai Zou, and Yu Cheng. You are catching my attention: Are vision transformers bad learners under backdoor attacks? In *ICCV*, 2023.
 - Yi Zeng, Won Park, Z Morley Mao, and Ruoxi Jia. Rethinking the backdoor attacks' triggers: A frequency perspective. In *ICCV*, 2021.
 - Mengxin Zheng, Qian Lou, and Lei Jiang. Trojvit: Trojan insertion in vision transformers. In *arXiv*, 2022.



⁷⁵⁶ B DETAILED SETTINGS FOR BACKDOOR DEFENSE

757 758 759

760

This section provides detailed information on the backdoor defenses applied in this paper. The settings of each defense are summarized as follows:

FT: We use AdamW (Loshchilov & Hutter, 2018) optimizer, the most popular optimizer for ViTs, to
fine-tune the backdoor ViTs for 20 epochs with a weight decay of 0.2. For ViT-B, the learning rate is
set as 3e-4. For other transformers, it is set as 5e-4. In addition, we adopt the cosine learning rate
schedule. Same as backdoor training, only simple data augmentations, including random crop with
padding and horizontal flipping, are used to retain the clean accuracy better and avoid the increasing
ASR of whole-image backdoor attacks caused by strong data augmentation as discussed in section 3.

FP: FP (Liu et al., 2018a) first prunes the last layer of CNNs by a predefined pruning threshold and then fine-tune the network on the clean subset of data. Similarly, we prune the last linear projection layer of transformer encoder blocks in ViTs. For the pruning partition threshold, we use *the tolerance of clean accuracy reduction* to limit the maximum drop of the benign accuracy following (Wu et al., 2022). In this paper, we set it to 0.9. The other settings are the same as the original paper (Liu et al., 2018a).

NAD: NAD (Li et al., 2021) first makes two copies of the original backdoor models, referred to as
the teacher model and student model respectively. Next, NAD fine-tunes the teacher model with the
vanilla FT. Finally, the finetuning of the student model is guided through neural attention transfer
from the teacher model. For the hyperparameter setting, we mainly keep in line with (Wu et al., 2022)
except for two differences: we train the student network for 20 epochs using the AdamW optimizer
instead of hundreds of epochs with SGD optimizer. The above changes are made because of the
observation shown in Appendix C and Appendix ??. As for the configuration of learning rate, we
follow FT, set 3e-4 for ViT-B and 5e-4 for other ViTs.

ANP: Wu et al. (Wu et al., 2020) observe that backdoor models are prone to output the target labels 781 when the neurons are perturbed by the adversarial perturbations. Inspired by this, they propose to 782 optimize the mask of each neuron, a continuous value in [0, 1], under adversarial neuron perturbations 783 and then prune neurons whose mask values are lower than the threshold, *i.e.*, hardening the continuous 784 mask values as binary masks. In this paper, we use the same settings as the original paper except 785 for applying 4000 iterations to avoid under-convergence of large models like ViTs (longer than the 786 2000 iterations for CNNs in the original paper). Compared to the hardened masks (pruned) applied in 787 their original paper, we find that soft masks, continuous mask values without hardening, can preserve 788 ACC better and decrease ASR further. Thus, we apply soft masks in this paper, and these masks are 789 applied to the channels of linear projection.

AWM: Compared to ANP, AWM (Chai & Chen, 2022) makes two improvements on CNNs. The authors apply soft element-wise weight masking instead of neuron pruning (hardened mask values) to avoid over-cutting beneficial information. Besides, they perturb the data instead of the neurons to utilize the training data more efficiently. When applied to ViTs, we mask the channel of the linear projection, similar to ANP. The other hyperparameters are the same as the original paper (Chai & Chen, 2022) without turning.

796 797

798

799 800

801

802

804 805 806

808 809

C THE EFFECT OF OPTIMIZER ON FP AND NAD

Table 8: The effect of optimizer on FP and NAD. AdamW gains higher ACC and lower ASR than SGD.

	((a) AC	С					(b) A	SR		
Attack		SC	GGD AdamW		Attack		SC	SGD		AdamW	
7 Hutek	No defense	FP	NAD	FP	NAD		No defense	FP	NAD	FP	NAD
Badnets	97.85	93.17	57.59	93.52	93.77	Badnets	100.00	0.90	4.24	0.91	1.57
Blend	97.85	93.41	94.27	92.59	94.09	Blend	100.00	9.67	48.57	0.73	8.94
CLB	97.83	27.20	94.31	93.22	93.88	CLB	96.23	8.21	10.15	1.70	7.27
SIG	97.50	77.34	94.31	93.88	93.86	SIG	90.57	1.93	5.00	0.81	3.60
AvgDrop	-	24.98	12.91	4.46↓	3.86↓	AvgDrop	-	91.53	79.71	95.66↑	91.36↑

In this section, we compare the performance of SGD and AdamW on the other two fine-tuning-based methods, FP and NAD, following the settings in section 3.2. As shown in Table 8, the results demonstrate that, compared to SGD, AdamW always performs better on FP and NAD. For example, SGD results in an average ACC drop of 24% in FP, much larger than 4.46% caused by AdamW.
Besides, SGD also has a little worse defense performance.

D THE ACCURACY OF OUR ATTACK ON CIFAR-10 DATASET

Table 9: ACC (%) of our attacks with different ViT variants on the benchmark dataset. The best results are in **bold**.

Defense	Attack			Vanilla						Ours		
Derense	7 Huek	ViT-B	DeiT-S	Swin-B	Cait-S	XciT-S	Vi	iT-B	DeiT-S	Swin-B	Cait-S	XciT-S
	BadNets	97.85	97.67	98.53	98.47	97.83	98	8.18	97.75	98.69	98.35	97.90
No defense	Blend	97.85	97.98	98.90	98.62	98.39	- 98	3.04	97.86	98.75	98.47	98.34
No defense	CLB	97.83	97.70	98.41	98.27	97.65	97	7.88	97.83	98.49	98.27	97.72
	SIG	97.50	97.44	98.56	98.21	98.05	93	7.88	97.36	98.67	98.14	97.89
	BadNets	93.79	93.17	95.23	94.24	93.99	94	4.03	93.47	95.17	95.17	93.89
ET	Blend	93.30	93.57	95.32	94.56	93.96	94	4.00	93.99	95.70	94.91	94.33
1.1	CLB	94.06	93.33	95.75	94.99	93.95	94	4.20	93.78	95.96	95.07	94.43
	SIG	93.51	93.75	95.47	94.75	94.07	93	3.38	93.51	95.10	94.97	94.44
	BadNets	93.52	93.40	95.84	95.18	94.57	93	3.67	93.41	95.98	95.29	93.59
ED	Blend	92.59	94.06	95.94	94.69	94.37	93	3.05	93.96	96.11	95.43	94.79
гг	CLB	93.22	93.99	95.91	95.36	94.55	93	3.15	94.17	95.48	95.42	94.36
	SIG	93.88	93.36	95.97	95.50	94.54	93	3.75	93.84	96.24	95.20	94.37
	BadNets	93.77	93.75	93.62	95.69	95.86	93	3.82	94.92	94.44	95.82	95.32
NAD	Blend	94.09	94.20	94.29	95.22	95.69	94	4.12	94.49	93.81	95.50	95.08
NAD	CLB	93.88	94.21	93.67	95.83	95.62	94	4.02	94.82	94.14	95.67	95.87
	SIG	93.86	94.30	93.90	95.91	95.90	93	3.95	94.08	94.75	95.55	95.15
	BadNets	94.26	93.77	98.05	96.89	93.45	94	4.40	94.49	97.62	97.39	95.78
AND	Blend	92.70	95.92	94.91	97.27	95.42	95	5.67	94.70	97.48	97.32	83.74
AINF	CLB	95.71	94.93	98.05	97.43	85.98	95	5.83	94.37	97.60	97.41	91.26
	SIG	92.60	95.25	97.79	97.76	96.18	94	4.62	94.70	97.51	96.98	96.55
	BadNets	95.02	94.52	96.39	95.93	95.46	93	3.87	94.91	96.28	96.18	95.43
AWM	Blend	95.08	94.99	93.00	96.51	96.00	95	5.06	94.82	95.38	96.28	94.40
AWM	CLB	95.60	94.94	95.20	96.17	95.33	95	5.12	94.84	94.22	96.41	95.53
	SIG	94.58	94.76	96.89	96.59	96.05	94	1.46	94.43	96.90	96.57	95.80

We have discussed the attack performance of our proposed method as shown in Table 4 of Section 5.1. Here, we continue to explore the effect on the accuracy of our attacks. As shown in Table 9, the backdoored models with our method have comparable accuracy to their baselines (without our method), which indicates our method does not influence the utility of the backdoored model and guarantees the stealthiness of backdoored models with our method.

E THE PERFORMANCES OF CAT AGAINST THE DETECTION-BASED DEFENSES

We further evaluate CAT on Neural Cleanse (NC) (Wang et al., 2019) to see whether CAT can help existing attacks better bypass the detected-based defense. NC is composed of two stages: Firstly, it reconstructs all possible triggers through optimization and determines whether the victim model is implanted with a backdoor via outlier detection. Secondly, if the answer is true, it will mitigate the backdoor behavior through unlearning with the reconstructed trigger, i.e., restoring the performance even with the presence of the trigger. We examine whether CAT can better bypass NC in these two stages, and all experiments are performed on the CIFAR-10 dataset with ViT-B and Swin-B architectures, covering both the white-box and black-box settings.

Detection Stage: NC reconstructs potential triggers for each class and uses the anomaly index metrics to determine if one of them is a valid trigger. The larger the anomaly index, the more likely it is to be a real backdoor trigger. Here, we calculate the anomaly indexes of the attack with or without CAT for comparison. The results in Table 10 (a) show that CAT can always lower anomaly indexes, making the attack stealthier. For example, the vanilla badnets attack obtains anomaly indexes of 7.45, which is larger than those after combining CAT (5.04). It means CAT can help existing attacks better bypass the detection of NC.

867	(a) Anon	(a) Anomaly Index			(b) ASR aft	(c) ACC after unlearning					
868		ViT-B	Swin-B			ViT-B	Swin-B			ViT-B	Swin-B
869 870	BadNets CAT+BadNets	7.45 5.04	4.17 3.25		BadNets CAT+BadNets	1.08 99.99	11.67 56.69		BadNets CAT+BadNets	96.85 97.22	96.87 96.35
871	Blend CAT+Blend	3.14 1.60	3.80 1.62		Blend CAT+Blend	0.66 53.49	1.24 21.88		Blend CAT+Blend	96.61 97.08	96.97 96.84
873	CLB CAT+CLB	7.13 2.48	2.99 2.26		CLB CAT+CLB	0.36 6.25	1.28 9.86		CLB CAT+CLB	96.78 96.75	96.88 96.90
874 875	SIG CAT+SIG	2.26 0.90	3.47 1.16		SIG CAT+SIG	5.64 43.79	2.36 15.63		SIG CAT+SIG	96.78 97.06	96.01 97.03

Table 10: Performance (%) of CAT against NC on the CIFAR-10 dataset.

876 877 878

879

880

882

883

884

885

886

887

889

899

900 901

864

> **Unlearning Stage:** Next, the defenders use the reconstructed triggers to mitigate the backdoor behavior once the reconstructed triggers are identified. Specifically, they fine-tune the model to predict ground-truth labels in the presence of the triggers, i.e., unlearning the backdoor behavior. Here, we explore whether CAT makes existing attacks more resistant to unlearning. According to previous research (Wu et al., 2022) which observes that the unlearning process of NC with CNNs' default settings will decrease the benign accuracy a lot ($\lambda 50\%$), we make the following adaptations based on the observations in our paper: (1) Use AdamW optimizer to unlearn the backdoored models. (2) Unlearn the backdoored model only for 20 epochs. We summarize the results in Table 10 (b). The table shows that CAT can make unlearning more difficult and keeps backdoor behavior inside the model. Therefore, we can conclude that CAT has a better capability of resisting the detection-based defense.

F THE SETTING OF OUR ATTACK ON IMAGENET DATASET



Figure 7: Examples for the benign and backdoor images on ImageNet dataset.

902 Attack: Since the huge computational cost, we fine-tune the pre-trained ViT-B on the poisoned 903 ImageNet with 512 batch size and 10 epochs to insert backdoors. Because ImageNet is a high-904 resolution dataset, we increase the trigger size of badnets attacks to 21×21 for better poisoning. 905 For the Blend attack, we resize the image of gaussian noise to 224×224 to accommodate the large 906 input size on ImageNet. In Figure 8, we show examples of benign and backdoor images. For other 907 settings of the vanilla poisoning, we keep the same with our experiments on CIFAR-10 (Please 908 refer to Appendix A for details.). For the settings of our proposed attack, we follow the settings of 909 CIFAR-10 except for the following two points: During the perturbation generation step, the budget and step size are set to 8/255 and 2/255, respectively. Similar to the vanilla backdoor attack, the 910 patch size of RMP is enlarged to 16 because ImageNet is a high-resolution dataset. For ViT-specific 911 attacks, we choose DeiT-B (Touvron et al., 2022) which has the exact same architecture as ViT-B for 912 poisoning without any hyperparameter change. 913

914 **Defense:** First, for the defense methods unrelated to architectures, to achieve a better acceleration 915 of the experiments on ImageNet, we adopt a large batch size of images for defense. In detail, for fine-tuning-based defense, the batch size is set to 512. For pruning-based defense, the batch size is 916 set to 128 to avoid the out-of-memory problem on 4 NVIDIA 3090 GPUs. Other settings are the 917 same as our experiment on CIFAR-10. Please refer to Appendix B for details. As for the ViT-specific

attack: attention blocking (AB), we adopt the default setting recommended by (Subramanya et al., 2024): during the inference stage, we block out the region of size 30×30 which is highlighted by Attention Rollout (Abnar & Zuidema, 2020).

G THE ACCURACY OF OUR ATTACK ON IMAGENET DATASET (THE WHITE-BOX SETTING)

Like the experiments on CIFAR-10, we also evaluate the effect of our method on ACC for large datasets like ImageNet. The results in Table 12 show that our method does not influence the utility of the backdoored models and the stealthiness of backdoored models on large datasets can also be further guaranteed.

Table 11: ACC (%)) of our attack on ImageNet data	set. The higher ACC is in bold .
-------------------	----------------------------------	---

Attack	Before	FT	FP	NAD	ANP	AWM	AB
TrojViT	80.59	76.82 7	6.93	77.55	76.31	77.78	-
DBIA	79.52	78.3	75.2	77.18	76.49	78.94	-
Badnets CAT+Badnets	80.82 81.01	71.05 6 71.41 6	58.10 5 8.31	72.38 72.69	69.56 69.79	76.40 76.62	74.86 74.51
Blend CAT+Blend	80.82 81.03	71.03 6 71.12 6	5 8.43 58.39	72.60 72.62	69.69 69.96	76.77 76.36	74.72 74.73

THE PERFORMANCE OF CAT ON THE IMAGENET DATASET UNDER THE Η **BLACK-BOX SETTING**

Table 12: Performance (%) of CAT on the ImageNet dataset under the black-box setting. The better performances are in **bold**. CAT obtains better performances than the vanilla attack.

	(a) ASR					(b) ACC							
Attack	Before	FT	FP	NAD	ANP	AWM	Attack	Before	FT	FP	NAD	ANP	AWM
Badnets	100.00	31.91	22.23	42.36	42.29	35.92	Badnets	83.02	77.28	76.06	77.45	68.93	75.78
CAT+Badnets	100.00	79.17	35.96	61.18	73.62	59.58	CAT+Badnets	83.11	77.27	76.10	77.86	69.54	76.84
Blend	100.00	6.10	3.01	18.00	10.85	31.96	Blend	82.93	76.70	76.37	77.81	68.23	76.22
CAT+Blend	100.00	21.62	22.32	35.11	43.84	45.92	CAT+Blend	83.09	76.87	76.51	77.26	70.16	75.96

ACTIVATION DIFFERENCE WITH OR WITHOUT COMBINING CAT Ι



