
Salient Conditional Diffusion for Defending Against Backdoor Attacks

Brandon B. May[†] N. Joseph Tatro^{†*} Dylan Walker Piyush Kumar Nathan Shnidman
Vision & Image Understanding
Systems & Technology Research
Woburn, MA 01801

Abstract

We propose a novel algorithm, **Salient Conditional Diffusion (Sanctifi)**, a state-of-the-art defense against backdoor attacks. **Sanctifi** uses a denoising diffusion probabilistic model (DDPM) to degrade an image with noise and then recover said image. Critically, we compute saliency-based masks to condition our diffusion, allowing for stronger diffusion on the most salient pixels. As a result, **Sanctifi** is highly effective at diffusing out triggers in data poisoned by backdoor attacks. At the same time, it reliably recovers salient features when applied to clean data. This performance is achieved without requiring access to the model parameters of the Trojan network, meaning **Sanctifi** operates as a black-box defense.

1 Introduction

With the increasing societal adoption of deep learning models, adversarial robustness, the ability of deep neural networks (DNNs) to withstand adversarial attacks, has quickly become a topic of interest in the machine learning community [Madry et al., 2017]. As the field develops, attention is being given to sophisticated attacks that align more closely to practical use cases. In this work, we focus on defending against backdoor attacks introduced in *BadNet* in Gu et al. [2017] as it is particularly challenging to defend against [Li et al., 2020]. In essence, a backdoor attack involves the passing of poisoned data, such as an image containing a visual trigger, to a malicious network that performs adversarially in the presence of said trigger. This creates a striking scenario where an adversary can hack with precision into a seemingly innocuous *Trojan network* that they have released to the public.

In this work, we: **(1)** Propose a novel defense against backdoor attacks, **Sanctifi**, that *purifies* input with a diffusion model (DDPM) conditioned on a mask derived from input dependent saliency maps. **(2)** Establish state-of-the-art performance among backdoor defenses. While **Sanctifi** is a black-box defense, needing no explicit access to the model parameters of the Trojan network, it is competitive with fine-pruning [Liu et al., 2018a] and Neural Attention Distillation [Li et al., 2021a]. **(3)** Demonstrate the utility of salient conditioning in our novel algorithm. We experimentally find that less salient parts of an image create a strong prior for the reverse diffusion process of a DDPM. This allows us to more reliably recover clean salient parts of an image.

Given its straightforwardness and ease-of-use, we are encouraged by the performance of salient conditional diffusion. In this work, we first review related work and provide background. Next, we motivate and formally introduce our algorithm **Sanctifi** and empirically validate its performance.

[†]Equal Contribution

*Corresponding author: joseph.tatro@str.us

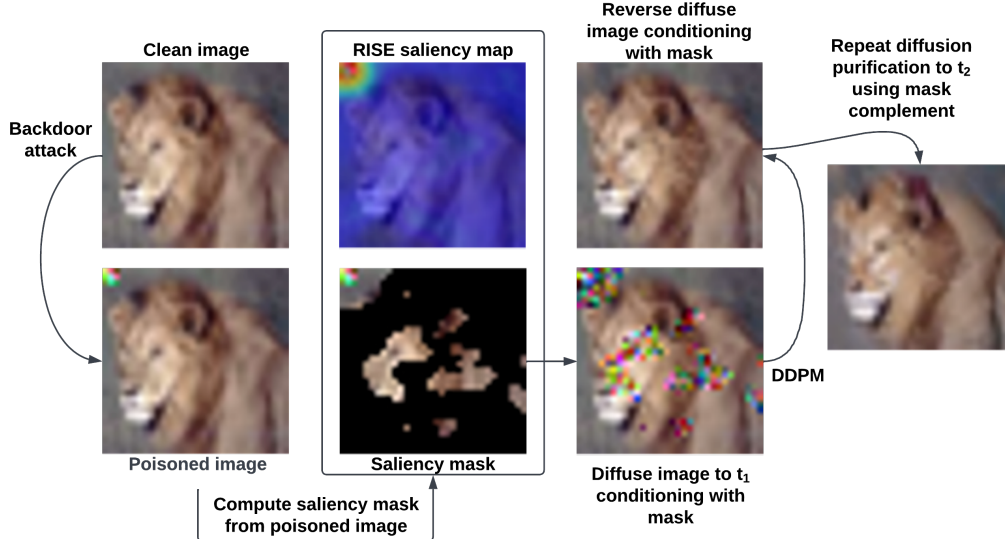


Figure 1: An illustration of Salient Conditional Diffusion (**Sanctifi**). First a trigger is added to a clean image in an attempted BadNet backdoor attack. With this input image, we compute saliency maps via RISE, and use the top-5 class maps to construct the binary visible mask, \mathbf{A} . Notice the trigger is unmasked as it is the most salient part of the top-1 map. We then apply diffusion purification, conditioned with the saliency mask, to the image for 300 time steps. Following this, we reapply diffusion purification using the reverse mask, $\mathbf{I} - \mathbf{A}$, with 100 time steps. Notice that **Sanctifi** is capable of diffusing out the backdoor trigger without largely degrading the entire image.

1.1 Related Work

Backdoor Attacks As stated previously, a backdoor attack involves an adversary training a malicious network such as BadNet in Gu et al. [2017]. Generally, this *Trojan network* contains a subnetwork that adversarially alters output in the presence of a specific trigger in the input. The detection of these Trojan networks is an ongoing research topic and is a difficult task given the vast permutations of possible subnetworks [Wang et al., 2020]. Recent work has demonstrated that even diffusion models are susceptible to backdoor attacks [Chou et al., 2022].

Surveyed in works including TrojanZoo [Pang et al., 2022] and BackdoorBox [Li et al., 2023], there are several types of defenses against backdoor attacks; input reformation, input filtering, model sanitization, and model inspection. Our proposed algorithm is an input reformation defense, meaning we prefilter network input while having no knowledge of model weights. Other state-of-the-art input reformation defenses include ShrinkPad (SP) [Li et al., 2021b]. Two state-of-the-art methods in backdoor defense, fine-pruning (FP) [Liu et al., 2018a], and Neural Attention Distillation (NAD) [Li et al., 2021a], are model sanitization defenses. These two methods are white-box methods that alter the model weights of the Trojan network. We stress that, as a black-box method, **Sanctifi** is applicable in more general real-world scenarios where fine-pruning is not. Concerning saliency-based methods, Februus [Doan et al., 2020] is a state-of-the-art method that uses a generative adversarial network (GAN) to inpaint an image after applying a GradCAM-derived mask. This use of Grad-CAM for trigger detection can also be found in SentiNet [Chou et al., 2020]. We will show later that **Sanctifi** successfully defends against more general attacks where Februus fails.

Diffusion Models and Conditioning Diffusion models such as denoising diffusion probabilistic models (DDPMs) [Ho et al., 2020] have quickly rivaled generative adversarial networks (GANs) [Goodfellow et al., 2020] in the task of image generation. DDPMs act by diffusing input through iteratively adding Gaussian noise and then learning the reverse diffusion process to recover the input image. This reverse diffusion process is able to generate data from noise in a Markov chain-like fashion. Research on diffusion models has exploded in the past few years. Kingma et al. [2021] analyzed the theoretical properties of the variational lower bound of DDPMs. Song et al. [2020] introduced an implicit variation, DDIM. Crossing into mainstream awareness, *Stable Diffusion* has powered generative AI apps used by the general public [Rombach et al., 2022].

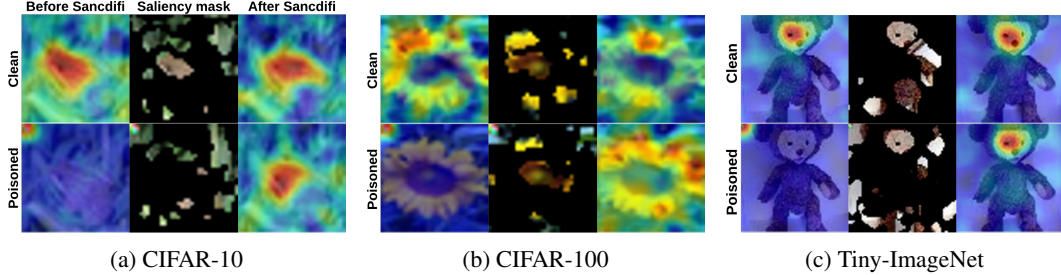


Figure 2: For CIFAR-10, CIFAR-100, and Tiny-ImageNet respectively, we show the computed saliency mask for ResNet-50. The top and bottom rows respond to clean data and BadNet attacked data. Each column displays; the top class saliency map, the computed saliency mask, and the top class saliency map after applying **Sancdifi**. Notice we have removed the saliency of the trigger.

Of particular interest, Nie et al. [2022] and Wu et al. [2022] both proposed diffusion purification, the use of diffusion models for adversarial purification against projected gradient descent (PGD) attacks. A strong inspiration for this work, DiffPure [Nie et al., 2022] is able to filter out PGD attacks with a small timescale diffusion and reverse diffusion process governed by a DDPM. We note that a defense appropriate for PGD attacks is not necessarily a valid defense against backdoor attacks [Weng et al., 2020]. Thus, we were inspired to see if diffusion purification could extend to backdoor attacks. In Section 3, we discuss the need for longer diffusion than in Nie et al. [2022] to sufficiently degrade triggers, while our novel salient conditioning is key to preventing a collapse in clean accuracy.

There has also been work that conditions the diffusion process of DDPMs. Dhariwal and Nichol [2021] introduced classifier-guided DDPMs. These models condition diffusion with a classifier gradient to steer the generative process to a user-specified class. Voynov et al. [2022] introduced a DDPM guided by the gradients of a latent edge predictor to improve text-to-image generation. These works focus on salient features when generating images. Unlike these methods, our work does not require user input, such as class, as a prior. Regarding mask-based conditioning, Aberman et al. [2022] introduced a model for reducing saliency within a region of an image determined via user-specified mask. Our algorithm has a different flavor as our mask is determined via saliency map and we do not decimate all salient features within the unmasked region of our images.

2 Salient Conditional Diffusion

To the best of our knowledge, this work is the first to propose the use of diffusion models (DDPMs) as a defense against backdoor attacks. A key contribution of **Sancdifi** is the use of saliency masks for conditioning diffusion purification. Before explaining our salient conditional diffusion algorithm, we state our backdoor attack model and discuss the motivation of our approach.

2.1 Backdoor Attack Model

In this work, we consider malicious networks that are sensitive to the presence of *triggers* in data. Such a backdoor trigger, a small 3x3 patch, can be seen in Figure 1 which illustrates the entire salient conditional diffusion process. As defined in Pang et al. [2022], the trigger starts with a pattern $p(\mathbf{x})$, that may depend on the data $\mathbf{x} \in \mathbb{R}^d$, of transparency value $\alpha \in [0, 1]$. Additionally, certain elements of the data are masked from the trigger following a mask template, $\mathbf{m} \in \{0, 1\}^d$. The trigger embedded data, $\mathbf{x} \oplus \mathbf{r}$, with trigger \mathbf{r} is defined as:

$$\mathbf{x} \oplus \mathbf{r} := (1 - \mathbf{m}) \odot ((1 - \alpha)\mathbf{x} + \alpha p(\mathbf{x})) + \mathbf{m} \odot \mathbf{x}. \quad (1)$$

Trojan networks are expected to handle both this poisoned data, $\mathbf{x} \oplus \mathbf{r}$, and clean data, \mathbf{x} . The clean data is associated with a label, y , while the target label for poisoned data is t . Let f_θ denote our Trojan network with parameters θ , and let \mathcal{L} denote the loss function. Then the objective of the Trojan model is to solve the following optimization problem,

$$\min_{\mathbf{r} \in \mathbb{R}, \theta} \mathcal{L}(\mathbf{x} \oplus \mathbf{r}, t) + \lambda \mathcal{L}(\mathbf{x}, y). \quad (2)$$

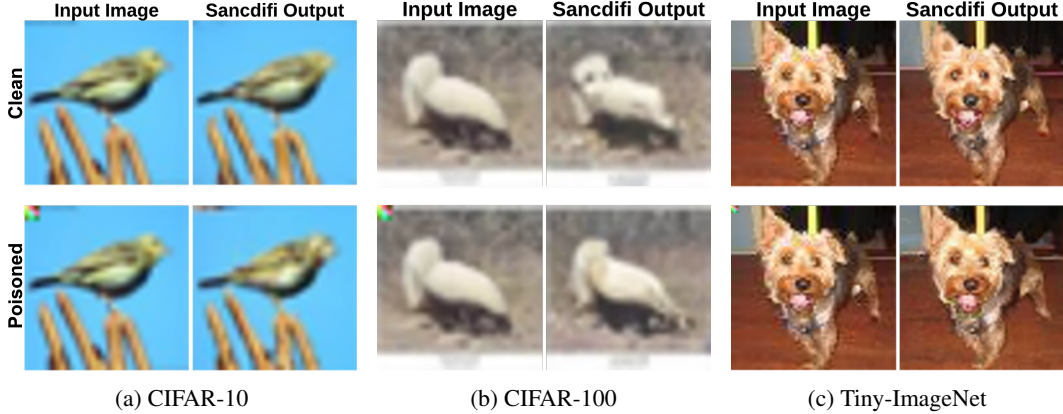


Figure 3: **Sancdifi** defending against BadNet attacks on CIFAR-10, CIFAR-100, and Tiny-ImageNet for ResNet-50. The top and bottom rows display it operating on clean and the corresponding BadNet attacked images respectively. Our method removes the trigger from poisoned data while preserving the integrity of clean data.

Here, R denotes the set of candidate triggers and λ is a hyperparameter managing the trade-off between clean accuracy and backdoor attack success.

We use BadNet [Gu et al., 2017], TrojanNN [Liu et al., 2018b], and WaNet [Nguyen and Tran, 2021] as our attack models. In BadNet, the trigger r is fixed, in contrast to other attacks that optimize some combination of the parameters in equation 1. For instance, TrojanNN optimizes the pixel values of the trigger to maximize certain neuron activations in the Trojan network. While the first two have input-invariant triggers, WaNet is input-adaptive and uses an imperceptible warping-based trigger. All attacks retrain a benign model f where a subset of the data has been poisoned. We also consider *Invisible BadNet*, an attack where the trigger has image-wide support but is L_∞ -bounded to make it visually imperceptible.

2.2 Motivation

As visualized in Figure 1, one of the most common trigger varieties in backdoor attacks are localized triggers. In the figure, the trigger is a 3x3 pixel signature within a 32x32 pixel image. Let $\hat{r} := x \oplus r - x$ denote the visible trigger. Then the support of \hat{r} is fairly concentrated. Due to time-frequency duality as described in Chaparro [2015], the frequency spectra of \hat{r} can be expected to have wide support. This has been confirmed empirically in Zeng et al. [2021].

It is well-known that (forward) diffusion operates on a function f by exponentially decaying its component frequency modes, with the rate of decay scaling with frequency. It follows diffusion can be used to degrade the high-frequency rich trigger associated with a backdoor attack. Unfortunately, diffusion also degrades the rest of the image. Then it is the reverse diffusion process of a DDPM, trained on clean data, that can recover the diffused image. Since images poisoned with this trigger are assumed not to be in the training data for the DDPM, we expect the trigger to be significantly degraded. It will be this degradation that prevents the activation of the the Trojan network. In order to minimize image degradation from the process, we would like to restrict diffusion to pixels associated with the Trojan trigger or those most likely to be recovered by the DDPM. **Sancdifi** is motivated by the idea that these pixels correspond with the most salient pixels in an image.

2.3 Methodology

With the use of DDPMs for defending against backdoor attacks motivated, we describe the **Sancdifi** algorithm, with a summary being available in Algorithm 1 and visible in Figure 1. A core component of our algorithm is the use of saliency to condition diffusion purification.

A saliency map S_k for a given image x , class k , and classifier network f measures the importance of each pixel of x . This importance is relative to f 's determination of the k -class probability of x . Arguably the most well-known saliency map algorithm is the white-box algorithm Grad-CAM, which

Algorithm 1 Salient Conditional Diffusion algorithm with image \mathbf{x} , Trojan network f , N RISE masks, time steps $\{T_1, T_2\}$, saliency percentile cutoff d , and r of top- r performance.

Require: $T_i \geq 0, d \in (0, 1), r \geq 1, i = 1$, trained DDPM to parameterize $p(\hat{\mathbf{x}}_{T-1}|\hat{\mathbf{x}}_T)$

$\mathbb{C} \leftarrow \text{top-}k(f(\mathbf{x}), r)$ indices

$\mathbb{S} \leftarrow \{\text{RISE}(\mathbf{x}, f, N, c), c \in \mathbb{C}\}$ ▷ See [Petsiuk et al., 2018] for RISE algorithm

$\mathbb{M} \leftarrow \{\mathbf{S}_i \leq \text{percentile}(\mathbb{S}_i, d), \mathbf{S}_i \in \mathbb{S}\}$

$\mathbf{A} \leftarrow \prod_i \mathbf{M}_i, \mathbf{M}_i \in \mathbb{M}$

while $i \leq 2$ **do**

$\mathbf{z} \leftarrow \text{sample } q(\mathbf{x}_{T_i}|\mathbf{x}_0)$ ▷ defined in equation 4

$\hat{\mathbf{x}}_{T_i} \leftarrow \mathbf{A}\mathbf{x}_0 + (\mathbf{I} - \mathbf{A})\mathbf{z}$

while $T_i \neq 0$ **do**

$\mathbf{z} \leftarrow \text{sample } p(\hat{\mathbf{x}}_{T_i-1}|\hat{\mathbf{x}}_{T_i})$ ▷ defined in equation 6

$\hat{\mathbf{x}}_{T_i-1} \leftarrow \mathbf{A}\hat{\mathbf{x}}_{T_i} + (\mathbf{I} - \mathbf{A})\mathbf{z}$

$T_i \leftarrow T_i - 1$

$\mathbf{A} \leftarrow \mathbf{I} - \mathbf{A}$

$i \leftarrow i + 1$

defines saliency as the gradient of the k -class classifier output $\nabla_{\mathbf{x}} f_k$ [Selvaraju et al., 2017]. We measure saliency using maps generated by the RISE algorithm [Petsiuk et al., 2018]. RISE saliency maps are computed in a black-box fashion, approximating Grad-CAM output, while requiring no knowledge of the model parameters of the network.

Given an input image \mathbf{x} , **Sanctifi** starts by computing the RISE saliency maps of \mathbf{x} for the top r classes determined by the Trojan network f_θ . Examples of RISE saliency maps can be seen in Figure 2. The most probable saliency map for clean images highlights meaningful pixels such as the body of a frog, the petals of a flower, or the face of a stuffed animal. In contrast, the most probable map for BadNet-poisoned images has the strongest response on the trigger. Encouragingly, we will show in Section 3.1 that the application of our defense to poisoned images produces saliency maps close to their clean counterparts.

From the k -class saliency map \mathbf{S}_k , we threshold the top d percentile of values to create a k -class saliency mask, \mathbf{M}_k . We desire our algorithm to have robust performance over different validation metrics such as top-5 accuracy. With that in mind, given the set of masks corresponding to the top- r most probable classes \mathbb{S}_M , we can define a composite saliency mask \mathbf{A} as their elementwise product. Concretely,

$$\mathbf{A} := \prod_{M \in \mathbb{S}_M} M \quad \text{where} \quad \mathbf{M}_k := \mathbf{S}_k \leq \text{percentile}(\mathbf{S}_k, d). \quad (3)$$

We will use \mathbf{A} to condition our diffusion processes. Intuitively, the composite mask ignores all but the most salient pixels of the most likely classes.

With the creation of our saliency mask, \mathbf{A} , we discuss the diffusion purification process motivated by Nie et al. [2022]. Our method of diffusion is taken from OpenAI’s improved-diffusion DDPM [Nichol and Dhariwal, 2021]. Given input data \mathbf{x} , we begin by diffusing it to time t by sampling from the distribution $q(\mathbf{x}_t|\mathbf{x}_0)$ and then applying the mask,

$$q(\mathbf{x}_t|\mathbf{x}_0) := \mathcal{N}(\mathbf{x}_t; \sqrt{\hat{\alpha}_t}\mathbf{x}_0, (1 - \hat{\alpha}_t)\mathbf{I}) \quad \text{where} \quad \hat{\alpha}_t := \prod_{i=0}^t (1 - \beta_i) \quad (4)$$

$$\hat{\mathbf{x}}_t = \mathbf{A}\mathbf{x}_0 + (1 - \mathbf{A})\mathbf{z}, \quad \mathbf{z} \sim q(\mathbf{x}_t|\mathbf{x}_0). \quad (5)$$

Here β_t denotes the variance schedule of the diffusion process. In practice, we define β_t to linearly increase over time. Note that the effect of applying the saliency mask to the output is equivalent to conditioning the diffusion process.

Table 1: **Sanctifi** (SD) defense results on BadNet for ResNet-50. Our metrics include clean accuracy reduction (CAR) and attack success rate (ASR) for top-1 and top-5 class performance. CAR represents the drop in accuracy for clean images after applying our defense. We desire low values of CAR and ASR. The other algorithms, referenced in Section 3, are fine-pruning (FP), ShrinkPad (SP), Neural Attention Distillation (NAD), and Februs (FB) respectively. Our algorithm has performance comparable to white-box defenses, FP and NAD, while having our CAR is more consistent across datasets compared to the other black-box defenses.

		Backdoor Defenses								
		top-1					top-5			
Dataset	Metric	SD	FP	SP	NAD	FB	SD	FP	SP	NAD
CIFAR-10	CAR	2.0	-1.0	1.0	13.0	13.0	0.0	0.0	0.0	1.0
	ASR	12.0	36.0	11.0	11.0	11.0	55.0	95.0	60.0	57.0
CIFAR-100	CAR	10.0	15.0	7.0	13.0	—	8.0	5.0	1.0	5.0
	ASR	1.0	1.0	1.0	1.0	—	8.0	3.0	12.0	5.0
Tiny ImageNet	CAR	7.0	0.0	16.0	10.0	—	5.0	0.0	6.0	7.0
	ASR	3.0	1.0	1.0	5.0	—	7.0	6.0	3.0	17.0

Table 2: **Sanctifi** (SD) defense top-1 accuracy on BadNet for other networks. Methods and metrics defined are in Table 1. We see that our performance extends to architectures beyond ResNet-50 and is comparable to other defenses. Particularly, we outperform ShrinkPad on these other architectures.

		Backdoor Defenses						
		CIFAR-100				CIFAR-10		
Network	Metric	SD	FP	SP	NAD	SD	FP	FB
Efficient-Net	CAR	13.0	18.0	0.0	4.0	4.0	-1.0	1.0
	ASR	1.0	1.0	35.0	21.0	10.0	11.0	10.0
ViT	CAR	17.0	1.0	3.0	—	5.0	0.0	1.0
	ASR	1.0	1.0	67.0	—	8.0	10.0	10.0

Letting $\hat{\beta}_t$ be $\frac{1-\hat{\alpha}_{t-1}}{1-\hat{\alpha}_t}\beta_t$, the reverse diffusion process involves both the following prior and posterior conditional distributions, p and q :

$$p(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu(\mathbf{x}_t, t), \hat{\beta}_t \mathbf{I}), \quad \mu := \frac{1}{\sqrt{1-\beta_t}} \left(\mathbf{x}_t - \frac{\beta_t}{1-\hat{\alpha}_t} \mathcal{E}_\theta(\mathbf{x}_t, t) \right), \quad (6)$$

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \hat{\mu}(\mathbf{x}_t, \mathbf{x}_0), \hat{\beta}_t \mathbf{I}), \quad \hat{\mu} := \frac{\sqrt{\hat{\alpha}_{t-1}}\beta_t}{1-\hat{\alpha}_t} \mathbf{x}_0 + \frac{\sqrt{1-\beta_t}(1-\hat{\alpha}_{t-1})}{1-\hat{\alpha}_t} \mathbf{x}_t. \quad (7)$$

To determine $\hat{\mathbf{x}}_0$, we iteratively sample from the prior distribution p . In the DDPM framework of Ho et al. [2020], the function \mathcal{E}_θ is parameterized by a neural network. Training a DDPM involves optimizing \mathcal{E}_θ to minimize the sum of KL-divergences of the conditional posteriors from the conditional priors at each time step in the reverse diffusion. Conditioning the reverse process is analogous to equation 5, with the critical difference being that masking must take place at each time step. Conditioning diffusion via a mask can be seen in Lugmayr et al. [2022] for further reference.

Following this, we reapply the diffusion purification on the resulting image using the complement of our salient mask, $\mathbf{I} - \mathbf{A}$, with time \hat{t} where $\hat{t} < t$. This safeguards against attacks with support across the entire image. We diffuse a shorter amount of time as we believe these low-saliency features are less recoverable by the DDPM. We show the importance of this second purification in Section 3.1.

Table 3: For transferability, **Sanctifi** results for ResNet-50 using the pretrained ImageNet DDPM. The CIFAR classes are subclasses of ImageNet classes. ImageNet is more complex than the CIFAR as it contains much more data. At the same time, images from ImageNet are in higher resolution. We find that there is some transferability in using the ImageNet DDPM to purify backdoor attacked data from the other related datasets. Both CAR and ASR are higher than in Table 1, but as the BadNet attack is successful 100% of the time in lieu of any defense, the transfer performance is notable.

Dataset	Metric	baseline top-1 (Table 1)	top-1	top-5
CIFAR-10	CAR	2.0	8.0	0.0
	ASR	12.0	17.0	59.0
CIFAR-100	CAR	10.0	24.0	14.0
	ASR	1.0	6.0	21.0

3 Numerical Experiments

Now that we have discussed the mechanics of our algorithm, we outline the setting of our numerical experiments to validate the performance of **Sanctifi**. For our experiments, we concern ourselves with the task of image classification in the presence of backdoor attacks. We employ multiple **datasets**; CIFAR-10, CIFAR-100, [Krizhevsky et al., 2009], and Tiny-ImageNet [Le and Yang, 2015]. Also, we experiment with multiple **architectures**; ResNet-50 [He et al., 2016], EfficientNet-B7 [Tan and Le, 2019], and a transformer ViT-Base-16 [Dosovitskiy et al., 2020]. Any metrics reported are for the dataset validation subsets. The networks cover a range of respective qualities; classical and useful, high-performing with low parameterization, and state-of-the-art and complex. Computationally, all experiments were performed on a NVIDIA RTX 2080 Ti GPU.

Concerning our algorithm, we use pretrained DDPMs from OpenAI’s improved-diffusion repository [Nichol and Dhariwal, 2021]. The DDPM we use for a given dataset is trained on that clean dataset. We rescale the timescale of the models from 4000 maximum steps to 1000 steps. For the CIFAR datasets, we diffuse out to 300 time steps for the first diffusion purification. For Tiny-ImageNet, we use 450 time steps as we find that it is needed to sufficiently defend against BadNet attacks. For the second diffusion purification step using the complement mask, $I - \mathbf{A}$, we diffuse out to 100 time steps. Our backdoor attacks are generated using the TrojanZoo suite [Pang et al., 2022] with their default parameters. In the case of WaNet, we use the BackdoorBox suite [Li et al., 2023].

Regarding saliency, we compute RISE maps using 2000 random binary masks. For the saliency threshold, we set a value of 95%. This cutoff is likely lower than necessary as the 3x3 trigger occupies less than 1% of an image in our experiments. The composite saliency map is aggregated across the top-5 classes to align us with top-5 metrics. For comparison, we evaluate four other defenses against the attacks mentioned in Section 1.1; fine-pruning (FP) [Liu et al., 2018a], ShrinkPad (SP) [Li et al., 2021b], Neural Attention Distillation (NAD) [Li et al., 2021a], and Februus (FB) [Doan et al., 2020]. Februus results are limited to CIFAR-10 due to availability of inpainting GANs.

3.1 Results

Performance on BadNet Attack We first discuss the results of our **Sanctifi** algorithm in defending against a BadNet attack on ResNet-50. Table 1 displays performance in terms of clean accuracy reduction (CAR) and attack success rate (ASR). While we focus mainly on top-1 classification accuracy, we also include top-5 classification results. To be clear, CAR denotes the reduction in accuracy on clean images after applying the defense algorithm. Intuitively, we desire CAR and ASR to be as low as possible. **Sanctifi** performs comparably to these state-of-the-art approaches. In this case, the winner among these methods is largely a question of the tradeoffs between CAR and ASR as well as top-1 and top-5 performance. The **Sanctifi** defense associated with these results is visible in Figure 3. In the figure, the BadNet trigger has clearly been diffused after purification. In contrast, the other salient regions of the images not covered by salient mask \mathbf{A} have been reliably recovered by the DDPM.

To further validate our performance, we repeat the previous experiment for our other architectures on CIFAR-100. We include a CIFAR-10 results for comparison with Februus . Summarized in Table 2,

Table 4: **Sancdifi** (SD) results on other backdoor attacks for ResNet-50. We also include PGD attacks. Clearly, our algorithm can handle other backdoor attacks such as the input-adaptive WaNet as well as the traditional PGD attack. So **Sancdifi** can be used for both backdoor and adversarial robustness. Our worst scenario is the image-wide Invisible BadNet attack, though we can resolve this by running the second diffusion for longer than 100 steps. For verification, see Table 7.

		Backdoor Defenses					
		CIFAR-100				CIFAR-10	
Attack	Metric	SD	FP	SP	NAD	SD	FB
WaNet	CAR	14.0	17.0	20.0	0.0	5.0	36.0
	ASR	2.0	18.0	38.0	97.0	1.0	74.0
Invisible BadNet	CAR	5.0	-3.0	1.0	3.0	3.0	2.0
	ASR	20.0	1.0	19.0	13.0	11.0	88.0
TrojanNN	CAR	9.0	4.0	8.0	2.0	5.0	8.0
	ASR	1.0	2.0	0.0	0.0	35.0	100.0
PGD	CAR	18.0	—	3.0	8.0	0.0	7.0
	ASR	0.0	—	4.0	95.0	10.0	88.0



Figure 4: Comparison of diffusion purification (**DiffPure**) to **Sancdifi**, which includes salient conditioning. This example is with CIFAR-100 and ViT. We display the clean image, the clean image purified with **Sancdifi**, the BadNet attacked image purified with **Sancdifi**, and the BadNet attacked image purified with DiffPure. Without salient conditioning, the face is destroyed.

the behavior of our algorithm is similar to Table 1. Notably, the performance of ShrinkPad does not generalize to other architectures. This shows that the performance of **Sancdifi** generalizes to various classes of neural networks.

Regarding the dataset-specific DDPM, Table 3 demonstrates that we still achieve decent performance when using the ImageNet DDPM across datasets. This suggests that a more general domain DDPM can be used in defending related, more specific data. This is helpful as in practice, we may wish to forgo training a DDPM from scratch if a pretrained one for similar data exists.

Performance against Other Attacks We have established that **Sancdifi** achieves competitive performance on BadNet backdoor attacks. For thoroughness, we consider other backdoor attacks; the input-adaptive WaNet, Invisible BadNet, and TrojanNN. We also include the traditional PGD adversarial attack. Table 4 contains our results on defending against these attacks. Our algorithm performs well across the various attacks. We can rectify our weakest performance, which is on Invisible BadNet, by increasing the complement diffusion purification steps past 100 iterations. Table 7 in the appendix verifies this is possible without notably harming our effectiveness against regular BadNet attacks. It is apparent that our defense can prevent PGD attacks better than adversarial retraining of Madry et al. [2017] which has CAR/ASR of 16.0%/8.0% for CIFAR-100. This is important as our algorithm is able to avoid the trade-off between adversarial robustness and backdoor robustness which has been suggested in literature [Weng et al., 2020]. Critically, Februous does not adequately defend against PGD attacks as its inpainting procedure only alters a small portion of the image, in contrast to **Sancdifi** which diffuses the entire image to some extent. Furthermore, inpainting

Table 5: Diffusion results **without** salient conditioning for ResNet-50. This reduces to the DiffPure algorithm [Nie et al., 2022]. Diffusion times are denoted relative to the maximum 1000 time steps. As diffusion time increases, ASR decreases at the cost of increased CAR. At less than 30% diffusion, ASR can become too high as in the case of Tiny-ImageNet. Yet the high diffusion leads to worse CAR. Notice that in the case of CIFAR-100, CAR is much higher at 30% than our algorithm (SD) in Table 1. Thus, saliency masking is needed.

Dataset	Metric	Diffusion Times					
		top-1			top-5		
		10%	20%	30%	10%	20%	30%
CIFAR-10	CAR	5.0	5.0	9.0	1.0	1.0	2.0
	ASR	90.0	14.0	11.0	100.0	63.0	61.0
CIFAR-100	CAR	13.0	31.0	47.0	2.0	15.0	31.0
	ASR	42.0	1.0	0.0	82.0	3.0	3.0
Tiny ImageNet	CAR	2.0	2.0	13.0	0.0	3.0	6.0
	ASR	99.0	47.0	9.0	99.0	53.0	15.0

Table 6: **Sanctifi** results for CIFAR-100 **without** second diffusion purification using the complement mask, $I - A$. We consider attack scenarios from Table 4. Invisible BadNet suffers most without the complement diffusion purification, with PGD also increasing in ASR.

Attack	Metric	baseline top-1	top-1	top-5
Invisible BadNet	CAR	5.0	1.0	6.0
	ASR	20.0	84.0	99.0
TrojanNN	CAR	9.0	4.0	5.0
	ASR	1.0	1.0	4.0
PGD	CAR	18.0	13.0	6.0
	ASR	0.0	8.0	33.0

the entire image is not possible. With this, **Sanctifi** has an advantage over the other input-reformation defenses in terms of architecture generalization and adversarial robustness.

Impact of Saliency Masks One might assume that vanilla diffusion purification a la DiffPure [Nie et al., 2022] is sufficient against backdoor attacks. Table 5 provides results on ResNet-50 where we have performed no salient thresholding and omit the second diffusion purification step. Strikingly, CAR is much worse without salient masking. Notably, DiffPure at 30% has the worst CAR across all defenses for the CIFAR-100 dataset. Additionally, Table 5 varies the choice of diffusion time; 10%, 20%, and 30% of maximum time. Lower diffusion times reduce CAR but at the expense of an increased ASR. Thus, salient conditioning is a critical part of our algorithm for successfully defending against backdoor attacks. A comparison of the output with and without salient conditioning is visible in Figure 4. We can see that while it is not the most salient, the masked part of the image offers a strong prior for the DDPM. This prior allows us to more reliably recover the unmasked part of the image excluding the backdoor trigger.

We also verify the importance of the second diffusion purification step with the complement mask in Table 6. We see that compared to Table 4, ASR is much higher in Invisible BadNet as well as PGD. Figure 5 in the appendix displays how the saliency mask evolves throughout applying **Sanctifi** to an image that has suffered an *Invisible BadNet* attack. Notice the second diffusion purification shifts focus back towards the left side of the house and the right side of its roof.

4 Conclusion

We have presented salient conditional diffusion, **Sanctifi**, a state-of-the-art defense against backdoor attacks. Our algorithm is intuitive with wide generalization over various datasets and network architectures. It is a black-box defense, requiring no knowledge of the parameters of a trojan model, with the performance of state-of-the-art defenses like fine-pruning. We have confirmed its performance against both the classic BadNet attack, imperceptible and input-adaptive WaNet, and TrojanNN attack. Additionally, we find that our algorithm can be used as a preprocessing step to improve the adversarial robustness of a system. Thus we avoid sacrificing adversarial robustness in the pursuit of backdoor robustness.

Salient conditioning has played a major role in allowing us to diffuse out backdoor triggers without massive degradation to other parts of an image. We believe conditional diffusion will play a strong role in the future in defending against backdoor attacks.

Acknowledgments and Disclosure of Funding

This work was supported by the DARPA AIE program, Geometries of Learning (HR00112290078).

References

- Kfir Aberman, Junfeng He, Yossi Gandelsman, Inbar Mosseri, David E. Jacobs, Kai Kohlhoff, Yael Pritch, and Michael Rubinstein. Deep saliency prior for reducing visual distraction. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun 2022. doi: 10.1109/cvpr52688.2022.01923. URL <http://dx.doi.org/10.1109/CVPR52688.2022.01923>.
- L Chaparro. Frequency analysis: the fourier transform (chapter 5). *Signals and Systems Using MATLAB*, pages 333–396, 2015.
- Edward Chou, Florian Tramer, and Giancarlo Pellegrino. Sentinet: Detecting localized universal attacks against deep learning systems. *2020 IEEE Security and Privacy Workshops (SPW)*, May 2020. doi: 10.1109/spw50608.2020.00025. URL <http://dx.doi.org/10.1109/SPW50608.2020.00025>.
- Sheng-Yen Chou, Pin-Yu Chen, and Tsung-Yi Ho. How to backdoor diffusion models?, 2022. URL <https://arxiv.org/abs/2212.05400>.
- Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis, 2021. URL <https://arxiv.org/abs/2105.05233>.
- Bao Gia Doan, Ehsan Abbasnejad, and Damith C. Ranasinghe. Februus: Input purification defense against trojan attacks on deep neural network systems. *Annual Computer Security Applications Conference*, Dec 2020. doi: 10.1145/3427228.3427264. URL <http://dx.doi.org/10.1145/3427228.3427264>.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale, 2020. URL <https://arxiv.org/abs/2010.11929>.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.
- Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Identifying vulnerabilities in the machine learning model supply chain, 2017. URL <https://arxiv.org/abs/1708.06733>.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun 2016. doi: 10.1109/cvpr.2016.90. URL <http://dx.doi.org/10.1109/cvpr.2016.90>.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.
- Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *Advances in neural information processing systems*, 34:21696–21707, 2021.
- Alex Krizhevsky et al. Learning multiple layers of features from tiny images. 2009. URL <https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf>.
- Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. 2015. URL http://vision.stanford.edu/teaching/cs231n/reports/2015/pdfs/yle_project.pdf.
- Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, and Xingjun Ma. Neural attention distillation: Erasing backdoor triggers from deep neural networks, 2021a.
- Yiming Li, Baoyuan Wu, Yong Jiang, Zhifeng Li, and Shutao Xia. Backdoor learning: A survey. *IEEE transactions on neural networks and learning systems*, PP, 2020.
- Yiming Li, Tongqing Zhai, Yong Jiang, Zhifeng Li, and Shu-Tao Xia. Backdoor attack in the physical world, 2021b.
- Yiming Li, Mengxi Ya, Yang Bai, Yong Jiang, and Shu-Tao Xia. Backdoorbox: A python toolbox for backdoor learning, 2023.

- Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. Fine-pruning: Defending against backdooring attacks on deep neural networks. *Lecture Notes in Computer Science*, page 273–294, 2018a. ISSN 1611-3349. doi: 10.1007/978-3-030-00470-5_13. URL http://dx.doi.org/10.1007/978-3-030-00470-5_13.
- Yingqi Liu, Shiqing Ma, Yousra Aafer, Wen-Chuan Lee, Juan Zhai, Weihang Wang, and X. Zhang. Trojaning attack on neural networks. In *Network and Distributed System Security Symposium*, 2018b.
- Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun 2022. doi: 10.1109/cvpr52688.2022.01117. URL <http://dx.doi.org/10.1109/CVPR52688.2022.01117>.
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- Anh Nguyen and Anh Tran. Wanet – imperceptible warping-based backdoor attack, 2021.
- Alex Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models, 2021. URL <https://arxiv.org/abs/2102.09672>.
- Weili Nie, Brandon Guo, Yujia Huang, Chaowei Xiao, Arash Vahdat, and Anima Anandkumar. Diffusion models for adversarial purification, 2022. URL <https://arxiv.org/abs/2205.07460>.
- Ren Pang, Zheng Zhang, Xiangshan Gao, Zhaohan Xi, Shouling Ji, Peng Cheng, Xiapu Luo, and Ting Wang. TrojanZoo: Towards unified, holistic, and practical evaluation of neural backdoors. *2022 IEEE 7th European Symposium on Security and Privacy*, Jun 2022. doi: 10.1109/eurosp53844.2022.00048. URL <http://dx.doi.org/10.1109/EuroSP53844.2022.00048>.
- Vitali Petsiuk, Abir Das, and Kate Saenko. Rise: Randomized input sampling for explanation of black-box models. *arXiv preprint arXiv:1806.07421*, 2018.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Bjorn Ommer. High-resolution image synthesis with latent diffusion models. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun 2022. doi: 10.1109/cvpr52688.2022.01042. URL <http://dx.doi.org/10.1109/CVPR52688.2022.01042>.
- 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*, pages 618–626, 2017.
- Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models, 2020. URL <https://arxiv.org/abs/2010.02502>.
- Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*, pages 6105–6114. PMLR, 2019.
- Andrey Voynov, Kfir Aberman, and Daniel Cohen-Or. Sketch-guided text-to-image diffusion models, 2022. URL <https://arxiv.org/abs/2211.13752>.
- Ren Wang, Gaoyuan Zhang, Sijia Liu, Pin-Yu Chen, Jinjun Xiong, and Meng Wang. Practical detection of trojan neural networks: Data-limited and data-free cases. In *European Conference on Computer Vision*, pages 222–238. Springer, 2020.
- Cheng-Hsin Weng, Yan-Ting Lee, and Shan-Hung Wu. On the trade-off between adversarial and backdoor robustness. In *Neural Information Processing Systems*, 2020.
- Quanlin Wu, Hang Ye, and Yuntian Gu. Guided diffusion model for adversarial purification from random noise, 2022. URL <https://arxiv.org/abs/2206.10875>.

Yi Zeng, Won Park, Z. Morley Mao, and Ruoxi Jia. Rethinking the backdoor attacks' triggers: A frequency perspective. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, Oct 2021. doi: 10.1109/iccv48922.2021.01616. URL <http://dx.doi.org/10.1109/ICCV48922.2021.01616>.

A Appendix

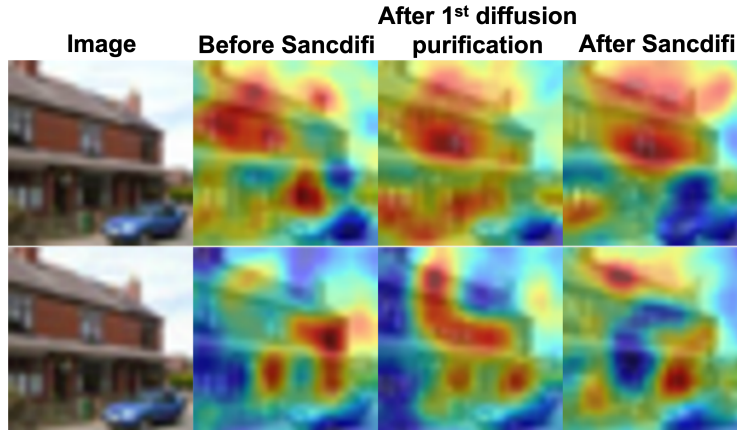


Figure 5: Saliency maps displayed for an **Invisible BadNet** attack on CIFAR100 for ResNet-50. The first and second rows correspond to the clean and attacked images. As the Invisible BadNet attack covers has support of the entire image, the first diffusion purification step does not sufficiently bring the saliency map inline with that of the clean image. We accomplish this with a second diffusion purification at smaller time using the complement mask.

Table 7: **Sancdifi** defense results on CIFAR-100 and ResNet-50 with different time steps used for the second diffusion purification. We find that a larger second timescale assists our algorithm in responding to Invisible BadNet attacks. At the same time, the increase in CAR on traditional BadNet attacks is still in line with other methods.

Diffusion times	Metric	Invisible BadNet		Traditional BadNet	
		top-1	top-5	top-1	top-5
300/100 steps	CAR	5.0	10.0	18.0	11.0
	ASR	20.0	72.0	0.0	7.0
300/150 steps	CAR	15.0	12.0	23.0	11.0
	ASR	3.0	33.0	0.0	7.0