

TnT Attacks! Universal Naturalistic Adversarial Patches Against Deep Neural Network Systems

Bao Gia Doan, Minhui Xue, Shiqing Ma, Ehsan Abbasnejad, Damith C. Ranasinghe

Abstract—Deep neural networks (DNNs), regardless of their impressive performance, are vulnerable to attacks from adversarial inputs and, more recently, Trojans to misguide or hijack the decision of the model. We expose the existence of an intriguing class of *spatially bounded, physically realizable, adversarial examples*—*Universal Naturalistic adversarial patches*—we call TnTs, by exploring the super set of the spatially bounded adversarial example space and the natural input space within generative adversarial networks. Now, an adversary can arm themselves with a patch that is naturalistic, less malicious-looking, physically realizable, highly effective—achieving high attack success rates, and universal. A TnT is *universal* because any input image captured with a TnT in the scene will: i) misguide a network (untargeted attack); or ii) force the network to make a malicious decision (targeted attack).

Interestingly, now, an adversarial patch attacker has the potential to exert a greater level of control—the ability to choose a location independent, natural-looking patch as a trigger in contrast to being constrained to noisy perturbations—an ability is thus far shown to be only possible with Trojan attack methods needing to interfere with the model building processes to embed a backdoor at the risk discovery; but, still realize a patch *deployable in the physical world*. Through extensive experiments on the *large-scale visual classification task*, ImageNet with evaluations across its *entire validation set* of 50,000 images, we demonstrate the realistic threat from TnTs and the robustness of the attack. We show a generalization of the attack to create patches achieving *higher attack success rates* than existing state-of-the-art methods. Our results show the generalizability of the attack to different visual classification tasks (CIFAR-10, GTSRB, PubFig) and multiple state-of-the-art deep neural networks such as WideResnet50, Inception-V3 and VGG-16. We demonstrate physical deployments in multiple videos at <https://TnTattacks.github.io/>.

I. INTRODUCTION

Deep Neural Network (DNN) systems are entrusted to make decisions in critical tasks such as self-driving cars [11], disease diagnosis [1] or face recognition [55], often, driven by their ability to achieve better-than-human performance. However, as DNN systems become more pervasive, it is creating the impetus for malevolent actors to attack them at: *i*) inference time—using *Adversarial Examples* such as unbounded perturbations [54], [25], [40], [41], [10] and adversarial patches (or spatially bounded perturbations) [6], [34]); and *ii*) training time—as in *Trojan attacks* [20], [23], [22]. The malintention of the adversary in each attack is the same—cause the DNN to fail by making an incorrect decision (untargeted attack) or

B. G. Doan, E. Abbasnejad, M. Xue and D. C. Ranasinghe are with the School of Computer Science, The University of Adelaide, Australia. Email: {giabao.doan; ehsan.abbasnejad; jason.xue; damith.ranasinghe}@adelaide.edu.au

S. Ma is with Department of Computer Science, Rutgers University, USA. Email: shiqing.ma@rutgers.edu

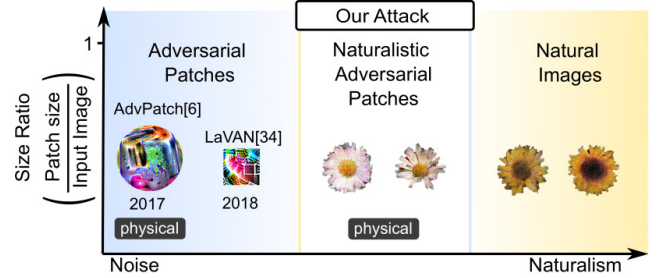


Fig. 1: An overview of the evolution of adversarial patch attack perturbation vectors. Our attack explores the hitherto elusive goal to realize visibly less malicious-looking adversarial patches—TnT—for: i) targeted attacks to misdirect any input to a target class; and ii) untargeted attacks.

manipulating the DNN to make the desired malicious decision (targeted attack).

Adversarial Examples are *unbounded input-specific* perturbations of additive *noise-like* vectors carefully crafted and applied to inputs to fool DNNs into making wrong classification decisions [25], [40], [10]. In contrast, unbounded **Universal Adversarial Perturbations (UAPs)** [41] or spatially bounded adversarial perturbations in LaVAN [34] have demonstrated the existence of *input-agnostic* additive noise patterns crafted to mislead the prediction of *any* input to a DNN. Notably, these *noise-like* additions from careful gradient perturbations are difficult to deploy in real-world settings, especially in perception systems. For example, variations in noisy physical environments can eliminate the noise vectors necessary to activate the DNN [39]. In contrast, the recent attack from Brown et al. [6] constructed a physically realizable **Adversarial Patch**, the AdvPatch in Fig. 1, to easily mount an *input-agnostic targeted* attack in the *physical world*; here, the attacker is constrained to perturbations *spatially bounded* to a region of an input to realize, a *printable* noise pattern.

On the other hand, an adversary mounting a **Trojan Attack** [20], [23], [59] *actively* poisons the training data or the network to embed a backdoor during the training process of a DNN. Unlike adversarial patches, the attacker relies on manipulating the construction of the network. Hence, the attacker is able to self-select *any* spatially bounded input of physical appearance and size—a *secret trigger* in Trojan vernacular or essentially a *patch*—to activate the backdoor, later, at the inference stage. Importantly, the trigger employed to misguide the backdoored classifier can now be *natural-looking* and *any object in a scene* from the natural image space shown in Fig. 1, such as a flower [20], can be .

An Interesting Observation. Despite the different methods

employed, an attacker armed with an adversarial patch or a Trojan trigger aims to cause the DNN system to fail—for example to misclassify an input or hijack the DNN predictions to achieve a desired target prediction. Notably, adversarial patches and Trojan triggers can misdirect a model in the presence of *any input class*—i.e. their effect is *universal* or input agnostic; however, unlike adversarial patches crafted from applying gradient perturbations to the input, a distinguishing facet of a Trojan attack highlighted by Bagdasaryan and Shmatikov [3] is the adversary’s *ability and freedom to self-select any secret trigger of naturalism, stealth, shape, size or features independently of the DNN model*.

*Remark 1: We seek to investigate the potential for a run-time attack with a **universal, physically realizable** patch allowing an adversary to exert a level of control, inconspicuousness and naturalism over the patch, thus far shown to be only possible with Trojan attack methods whilst **obviating** the need to interfere with the model building process and risk of discovery (Fig. 1)*

Such an attack would: i) *bridge* the divide between Trojan Attacks and Adversarial Attacks in the *input space*; and ii) constitute a pragmatic and inconspicuous zero-day exploit against already deployed deep perception models.

A. Our Attack Focus

Our investigation focuses on deep perception models and seek to further explore and understand new vulnerabilities. In particular, we seek to answer the following primary research questions (**RQ**) through our investigations:

RQ1: How can we discover TnTs that are physically realizable and naturalistic? (Section IV)

RQ2: How vulnerable are deep neural networks and their defended counterparts to TnT attacks? (Section VI & X)

RQ3: Do these TnTs have the features of *universality* or *input-agnostic nature* to misclassify any input to a targeted class? (Section VI-A)

RQ4: How *robust* are TnTs? (Sections VI-B & VI-D)

RQ5: How *generalizable* are TnT to unseen data or *transferable* to other networks? (Sections VI-A & VI-C)

RQ6: Can the effect of a TnT be explained by the occlusion caused by the patch or a network bias? (Section VI-E)

RQ7: What are the impacts of relaxing the need for naturalistic patches? (Section VIII)

RQ8: How comparable are patches generated from our attack method to state-of-the-art adversarial patch attacks? (Sections VIII, VIII-B & X)

RQ9: How robust are TnTs in the physical world? (Section IX)

B. Our Contributions and Results

This paper presents the results of our efforts to investigate generating adversarial patches that are less clearly malicious. We summarize our results and contributions as follows:

1) **We propose a new attack against DNNs.** Our attack method generates *Universal Naturalistic* adversarial patches, we call TnTs, by exploring the super set of the spatially bounded adversarial example space and the natural input space within generative adversarial networks (GANs). The TnTs we generate for attacking a DNN are:

- **Universal and Naturalistic.** A TnT is *universal* as any input with a TnT will fool the classifier and naturalistic, as assessed by a large cohort user study.
- **Highly effective in targeted and untargeted attacks against state-of-the-art DNNs.** In extensive experiments with ImageNet—a significant large-scale dataset with a million high-resolution images used for pre-training models of many real-world computer vision tasks—we achieved attack success rates of over 95% in the challenging attack setting of misclassifying *any* input to a *targeted* class.
- **Robust.** We observe high attack success rates irrespective of the location of the TnT; even with the TnT in a corner, i.e the *image background*.
- **Deployable in the physical world.** We conduct *physical world deployments* of our attack to demonstrate the practicability of the attack in various real-world settings. *Multiple, detailed, demonstration videos* are available at: <https://TnTattacks.github.io/>.
- **Highly generalizable.** A TnT discovered from 100 random sample images can effectively misguide the *entire* ImageNet validation set of 50,000 images. Further, we demonstrate effective attacks across multiple state-of-the-art networks (such as VGG-16, WideResNet50, SqueezeNet, ResNet18, MnasNet) and across 3 additional tasks: Face recognition (PubFig); Scene classification (CIFAR10); and Traffic sign recognition (GTSRB).
- **Transferable to mount black-box attacks.** We investigate attack transferability using the ImageNet classification task. We show that TnTs are *transferable* to other unknown network architectures for the same task (an attack in a *black-box* setting).
- **Highly effective at evading existing countermeasures against adversarial patch attacks.** We evaluate against both certifiable and empirical defenses.

2) **Our attack generalizes to generate physically realizable adversarial patches achieving higher attack success rates than state-of-the-art attacks.** When an attacker does not need naturalistic features, our attack leads to a new algorithm to generate *adversarial patches* of only 2% of the input image size with *higher* attack success rates; achieving a large margin of up to 44% compared to state-of-the-art adversarial patch attacks.

3) We demonstrate physical deployments in *multiple videos* at <https://TnTattacks.github.io/> and we contribute to the discipline by releasing the pre-trained networks and TnT artifacts to encourage future research and defenses.

II. BACKGROUND

Our attack formulation is based on employing generative adversarial networks (GANs) [26], therefore, we provide a brief background on GANs to aid with explaining our attack method in Section III, with details in Section IV.

Generative Adversarial Networks (GANs) have shown significant success in various applications of realistic image synthesis. A Generative Adversarial Network (GAN) [26] consists of a Generator G and a Discriminator D . Then if: i) \mathbf{x} is an input, which is an image of 3D tensor (width \times height \times depth); and ii) \mathbf{z} is a latent vector of dimension N , which is sampled from a noise distribution $P(\mathbf{z})$, in a classical GAN, the Generator G maps a source of noise $\mathbf{z} \sim P(\mathbf{z})$ to generate a synthetic image¹ $\tilde{\mathbf{x}} = G(\mathbf{z})$, while the Discriminator function is to discriminate the fake synthetic images $\tilde{\mathbf{x}}$ and the real ones \mathbf{x} , and the feedback from this Discriminator is utilized to help the Generator improve the image quality.

There are different methods to train a GAN, in this paper, we applied the Wasserstein GAN with Gradient Penalty (WGAN-GP) [27] since it has been shown to stabilize the GAN training process and improve the fidelity of samples. It involves solving the following optimization problem:

$$\min_G \max_D \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathbb{P}_g} [D(\tilde{\mathbf{x}})] - \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_r} [D(\mathbf{x})] + \lambda \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathbb{P}_{\tilde{\mathbf{x}}}} [(\|\nabla_{\mathbf{x}} D(\hat{\mathbf{x}})\|_2 - 1)^2],$$

where \mathbb{P}_r is the distribution of real images, and $\mathbb{P}_{\tilde{\mathbf{x}}}$ is the distribution of the interpolation between real and synthesized images. Here, \mathbb{P}_g is the conditional distribution of the synthesized images which we sample from the generator, that is, $\tilde{\mathbf{x}} = G(\mathbf{z})$, $\mathbf{z} \sim p(\mathbf{z})$. Using this min-max optimization, we learn \mathbb{P}_g to match \mathbb{P}_r . Samples from \mathbb{P}_r are drawn from a dataset of real objects; consequently, by sampling from the realized generator, we are able to obtain naturalistic image samples. We will exploit this ability of a GAN to synthesize naturalistic patches we call TnTs. We evaluate the naturalism of the generated TnTs, later, in Section VII.

III. OVERVIEW OF OUR ATTACK APPROACH

In this section, we provide an overview of our approach to provide a conceptual understanding the TnT attack method against Deep Neural Networks (DNNs) at the inference stage. First, we introduce the attack model, followed by our hypothesis, then our attack approach.

A. Attack Model

Our attacker strikes at the *inference* time, *i.e.* the attacker *does not intervene in the training process* in contrast to a Trojan attack. Consequently, the attacker does not leave any trace of tampering with the network to be discovered, making it relatively easy to deploy. Depending on the attacker’s knowledge of the target model to attack (*i.e.* neural architecture, inputs and outputs), we can consider either a *white-box* attack [25] where everything about the model is known, or *black-box* [45], [44] where *nothing* about the model is

¹The parameters are omitted for brevity but both generator and discriminator have their individual set of parameters.

known. In both cases, typically it is assumed that the attackers have access to some labeled training and validation data. The attacker also has access to computational resources to verify the method, *e.g.* a GPU-based cloud service.

White-box attacker. Following the attack models from prior adversarial patch attacks AdvPatch [6] and LaVAN [34], we primarily consider a *white-box* attack where the attackers have full access to the attacked network. Even if access to the model is not possible, or the model is not publicly available, one can employ a reverse engineering approach such as [57], [50], [8] to extract the model. Since defending against such attacks is challenging, it is of particular interest.

Black-box attacker. We also evaluate whether TnT can be exploited in the less restricted threat model of a black-box attack when the attacked network is unknown. We assume such an attacker has access to TnT obtained for the *same classification task* on a *different arbitrary network*. This allows an attacker to transfer the knowledge gained from, for example, a white-box attack on an arbitrary model using publicly available training data, to be used to attack a different model.

Goals. The attacker goals are to **(i)** exploit the vulnerability of a DNN to TnTs to extract TnT instances with **(ii)** high attack success rate (ASR), while **(iii)** maintaining the universality of the patch; the challenge is to discover a naturalistic patch.

B. Our Hypothesis

The decision boundaries learned by DNNs are subjected to the training examples presented to the network during the learning process. Consequently, DNNs must eventually approximate true decision boundaries learned from the training data. Unfortunately, due to the complex and highly non-linear structures within deep neural networks, it is hard to fully understand the learned boundaries within networks. For instance, Hendrycks et al. [30] reminds us that natural images captured from a scene with objects of classes within the training examples can sometimes cause errors and lead to incorrect predictions, even on “superhuman” ImageNet classifiers [29] despite having *no adversarial alterations*. Therefore, we can expect the existence of *adversarial patches* that are naturalistic but has the adversarial effect to alter the decision of the classifier. However, searching in the infinite space of all natural-looking small image patches is not feasible, therefore we constrain our search to the manifold of a Generative Adversarial Network (GAN) by taking inspiration from recent developments in GANs showing a tremendous ability to learn to generate realistic images [26].

*Remark 2: We hypothesize the existence of an intersection of the spatially bounded adversarial example space (*i.e.* patches) and the natural input space within generative adversarial networks (GANs) as illustrated in Figure 2.*

C. An Overview of Our Attack Approach

Since GANs are designed to map from a known (latent) distribution to the distribution of real images using a generator, as illustrated in Fig. 2, we consider the latent space \mathbf{z} of

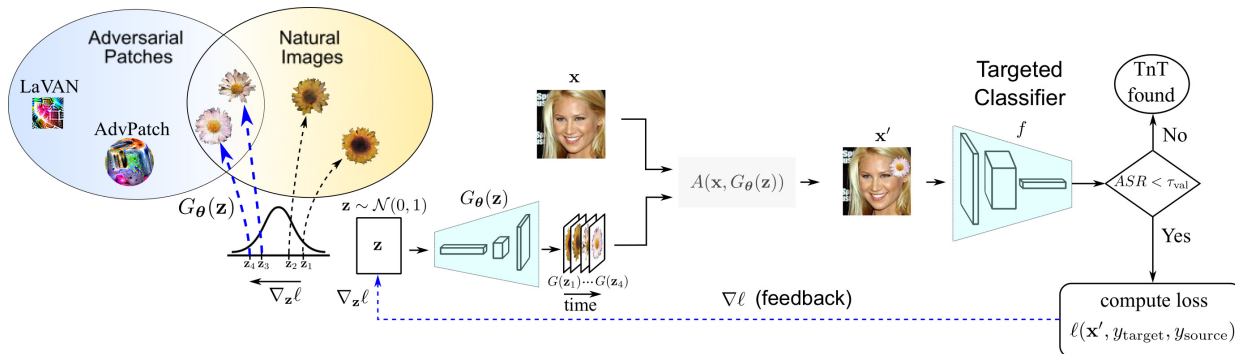


Fig. 2: Attack method for generating TnTs. Here, A is the patch stamping process, y_{target} is the targeted class designated by the attacker, y_{source} is the ground-truth label, $\ell(\mathbf{x}', y_{\text{target}}, y_{\text{source}})$ is the combined cross-entropy loss between the predicted score from the classifier f and the targeted as well as source label, and $\nabla \ell$ is the feedback from the Targeted Classifier f . The method is designed to iteratively approach high attack success TnTs by traversing through the latent space of the generator using gradient feedback.

the generator from which images are produced—as $G_{\theta}(\mathbf{z})$ —instead of searching in the infinite space of all natural-looking image patches; i.e. a standard Gaussian distribution $\mathcal{N}(0, 1)$ and has a much lower dimension in which traversal is easier. By getting feedback from the downstream classifier ($\nabla \ell$) we can **navigate** the **latent space** following the gradient feedback to *seek the latent vector from which a potential TnT can be generated*. Although not for the same attack, we acknowledge and discuss GAN-based attack approaches in Section XI.

Importantly, the learning algorithm we employ *determines* the best latent vector \mathbf{z} from which to generate a patch, a potential TnT because this latent space \mathbf{z} can capture the intrinsic structure of natural images from a simple latent distribution. Notably, since the generator was trained on natural images, samples from the latent space map to natural-looking image instances using a deterministic transformation (Generator). We demonstrate this process as an effective method to discover inputs capable of fooling the classifier whilst maintaining the naturalism of the patch. We will further support this claim by a large cohort user study, later, in Section VII. As illustrated in Fig. 2, we refer to this process as the *TnT Generator*. Effectively, our attack method takes advantage of a GAN’s ability to capture a rich distribution of natural images to discover the hypothesized region of the input space.

We distinguish our TnT attack from other Adversarial Patch attacks as follows.

- 1) Generated patches can be natural-looking and look less malicious than noisy perturbation patches in LaVAN [34] and AdvPatch [6].
- 2) Instead of *directly* applying perturbations to the input space that leads to noisy adversarial patches [6], [34], [38], we propose solving the problem of *searching for naturalistic patches* with adversarial effects by *indirectly* manipulating the *latent space* \mathbf{z} of a Generative Adversarial Network that has learnt to approach the natural patch distribution.
- 3) Different from the PS-GAN [38] attack based on *input-dependent* adversarial patches that *occluded* the salient features to realize an *untargeted* attack, our attack is:

- i) capable of both *targeted* and *untargeted* attacks;
- ii) *input-agnostic* (universal); and
- iii) *robust*—attack success is largely invariant to location, even at the corners or background of an image.

- 4) Our attack method generalizes to produce small, adversarial patches with *noise-like* additions of high attack success rates than such existing state-of-the-art attack methods with a large margin of up to 44%.
- 5) To the best of our knowledge, our study is the first to demonstrate an adversarial attack with a *universal*, *physically naturalistic* and *location independent* patch for *targeted* attacks in image classification tasks.

IV. TnT GENERATOR

In this section, we detail our TnT Generator method illustrated in Figure 2 for attacking a DNN with a concrete patch example. Without loss of generality, we propose using images of flowers to discover TnTs. Our primary motivation is that flowers exhibit synergy with a wide range of imaging scenarios and are unsuspecting, and therefore, inconspicuous. For example, someone might wear a flower T-shirt, wear a flower in their hair or hat to impersonate someone else in a face recognition system. Similarly, sticking an inconspicuous, natural-looking sticker on a traffic sign identical to [20] can fool a self-driving car to misclassify, for example, a STOP sign as a 100 mph Speed Limit sign with catastrophic consequences.

A. Training the Generator

An advantage of a GAN is that they are *unsupervised* models that only need unlabeled data that can be *cheaply* obtained. In our study, we collect a random, unlabelled flower image dataset from open-source Google Images [32] to build a flower dataset to learn the natural flower distribution. We selected the WGAN-GP loss function as it has shown to stabilize the training process of a GAN.

B. Transformation to a TnT Generator

To realize the TnT (representing a flower patch in our attack scenario), we need to *search* for flowers from our synthesized

flower distribution that has an adversarial effect on the attacked network. Here, we hypothesize that GANs have learned the super set of both natural-looking images and adversarials as illustrated in Fig. 2; therefore, by searching through this synthesized distribution, we expect to find a *structured, natural-looking* perturbation rather than a random noisy one. *First*, we formalize the notation of a TnT and, *second*, we propose a method for realizing such a TnT. Consider:

- y_{source} is the source class labeled for a given image \mathbf{x} .
- y_{target} is the targeted class designated by the attacker.
- The loss is between prediction and the label $\ell(\mathbf{x}, y_{\text{source}})$ for the *untargeted* and $\ell(\mathbf{x}, y_{\text{target}})$ for the *targeted* attack—cross-entropy loss function of the neural network, given image \mathbf{x} and a class label y . Notably, we intentionally omit network weights (θ) and other parameters in the cost function because we assume them to be fixed and remain untouched post network training.

Now, more formally, the attacker uses a trained model M that predicts class membership probabilities $p_M(y|\mathbf{x})$ to input images $\mathbf{x} \in \mathbb{R}^{w \times h \times c}$. We denote by $\mathbf{y} = p_M(\mathbf{x})$ the vector of all class probabilities, and by $y_{\text{argmax}} = \arg \max_{y'} p_M(y = y'|\mathbf{x})$ the highest scoring class for input \mathbf{x} (the classifier’s prediction on the *source class*). The attacker seeks an image \mathbf{x}' that fools the network such that $y' \neq y_{\text{source}}$ or $y' = y_{\text{target}}$ for $y' = \arg \max_{y'} p_M(y = y'|\mathbf{x}')$. The image \mathbf{x}' is composed of the original image with a *natural* patch stamped on it. We denote this process by a function $A(\mathbf{x}, G(\mathbf{z}))$.

In our TnT generation process, we firstly sample a vector $\mathbf{z} \sim p(\mathbf{z})$ with $\mathbf{z} \in \mathbb{R}^N$ where $N = 128$. This latent vector will be fed into our Generator pre-trained from Sec. IV-A in order to produce the flower patch $\delta = G(\mathbf{z})$ where $G: \mathbb{R}^N \rightarrow \mathbb{R}^{w \times h \times c}$. The flower patch is subsequently stamped at the lower-right corner of the image to avoid occluding the main features in alignment with the intentions in previous works [34], [49]. Later, in Section VI-B, we also evaluate the efficacy of the flower patch at nine different random locations. The size of the patch (flower) can be determined as necessary by the adversary to achieve their objectives (related to the attack success rate and TnT size discussed later in Section VIII-B). Based on the predefined location and patch size, we then utilize the *image thresholding* method of OpenCV [42] to determine the binary mask \mathbf{m} with $\mathbf{m}_{i,j} \in \{0, 1\}$ for i th row and j th column pixel of an image, to remove the background of δ . This patch is then affixed to the input image to obtain the adversarial sample \mathbf{x}' , i.e.:

$$\mathbf{x}' = (1 - \mathbf{m}) \odot \mathbf{x} + \mathbf{m} \odot \delta, \quad (1)$$

where \odot is the element-wise product.

To determine the ability of \mathbf{x}' to act as an adversary and receive feedback to choose a better latent variable, we test it with the trained neural network classifier. The sample \mathbf{x}' is fed into the classifier to obtain prediction scores for each individual class. The loss obtained from comparing the prediction scores and the target labels y_{target} , $\ell(\mathbf{x}', y_{\text{target}})$ as well as the source labels y_{source} , $\ell(\mathbf{x}', y_{\text{source}})$ are then calculated (e.g. using cross entropy). We use the additional loss $\ell(\mathbf{x}', y_{\text{source}})$ as it was shown to help the targeted attack converge faster. We

Algorithm 1: TnT Generator

```

1 Inputs: a batch of images  $\{\mathbf{x}^{(i)}\}_{i=1}^m$  with batch size  $m$ ,
   targeted label  $\{y_{\text{target}}^{(i)}\}_{i=1}^m$ , source label  $\{y_{\text{source}}^{(i)}\}_{i=1}^m$ ,
   model  $p_M$ , number of iteration  $n_{\text{iter}}$ , the learning rate
   parameter  $\epsilon$  to update the latent vector, the
   hyper-parameter  $\lambda$  to balance the loss, and the
   thresholds to detect the TnT  $\tau_{\text{batch}}$ ,  $\tau_{\text{val}}$  for batch and
   validation set respectively;
2 Initialization: fool = 0, detect = False;
3 while detect = False do
4   Sample a batch of images  $\{\mathbf{x}^{(i)}\}_{i=1}^m$ ;
5   Sample a latent variable  $\mathbf{z} \sim p(\mathbf{z})$ ;
6   for  $j = 1, \dots, n_{\text{iter}}$  do
7      $\delta = G_{\theta}(\mathbf{z})$   $\triangleright$  a flower patch;
8     Generate the mask  $\mathbf{m}$  based on  $\delta$ ;
9      $\delta' \leftarrow \text{bgremoval}(\delta, \mathbf{m})$   $\triangleright$  Background removal;
10    for  $i = 1, \dots, m$  do
11       $\mathbf{x}'^{(i)} = (1 - \mathbf{m}) \odot \mathbf{x}^{(i)} + \mathbf{m} \odot \delta'$ ;
12       $y_{\text{argmax}}^{(i)} = \arg \max_y p_M(y|\mathbf{x}'^{(i)})$ ;
13      if  $y_{\text{argmax}}^{(i)} = y_{\text{target}}^{(i)}$  then
14        | fool = fool + 1;
15      end if
16    end for
17  end for
18   $L = \ell(\{\mathbf{x}'^{(i)}\}_{i=1}^m, \{y_{\text{target}}^{(i)}\}_{i=1}^m) -$ 
19   $\lambda \ell(\{\mathbf{x}'^{(i)}\}_{i=1}^m, \{y_{\text{source}}^{(i)}\}_{i=1}^m)$ ;
20   $\nabla_{\ell} = \frac{\partial L}{\partial \mathbf{z}}$ ;
21   $\mathbf{z} = \mathbf{z} - \epsilon \text{sign}(\nabla_{\ell})$ ;
22  if fool >  $\tau_{\text{batch}}$  then
23    # further verify the realized TnT;
24    ASR  $\leftarrow$  Validate( $\delta, \mathcal{X}_{\text{val}}$ )  $\triangleright$  verify on  $\mathcal{X}_{\text{val}}$ ;
25    if ASR  $\geq \tau_{\text{val}}$  then
26      | detect = True;
27    end if
28  end while

```

then compute the gradient of this loss with respect to the latent variable \mathbf{z} , i.e. $\nabla_{\mathbf{z}} \ell(\mathbf{x}', y_{\text{target}}, y_{\text{source}})$. We then update the latent variable \mathbf{z} using gradient descent in the direction of minimizing this loss. Note that this does not change the classifier or GAN parameters and only updates the latent variable to increase the likelihood of attack success.

During TnT generating, for a random set of inputs, if a *threshold* percentage of inputs \mathbf{x}' can fool the network, the TnT is considered *universal*. The reason for setting a threshold here is to improve the algorithm’s speed, so that if the attack is successful in a batch, then we test whether it generalizes to the validation set \mathcal{X}_{val} . The complete process is summarized in the **Algorithm 1**.

V. ATTACK EXPERIMENT SETTINGS

We conduct an extensive experimental evaluation regime to evaluate our attack method. We describe the: i) Datasets; ii) GAN employed and training; iii) Attack configurations; and iv) Implementation Details employed in our quantitative evaluations in Section VI.

Datasets and Model Architectures. We employ popular real-world visual classification tasks. We conduct extensive experiments with the large-scale visual recognition dataset, ImageNet. Notably, the dataset is widely used as a pre-training model to achieve “superhuman” performance on classification tasks [29]. Additionally, to demonstrate the generalization of our method, we also evaluate on 3 other visual classification tasks: i) Scene Classification (CIFAR10); ii) Traffic Sign Recognition (GTSRB); and iii) Face Recognition (PubFig). Model architectures and testing samples used for each task are summarized in Table VIII. Model and dataset details are in **Appendix C**.

GAN Configuration and Training. To illustrate the significant threat posed, we demonstrate our attack method is easy to mount, low cost and does not require access to costly labeled data, and an attacker does not require specialized expertise in machine learning. Consequently: i) we utilized the off-the-shelf GAN framework of Pytorch, TorchGAN [43]; ii) used an off-the-shelf web crawler to automatically curate a dataset of random flower species from Google images; and iii) used only 945 flower images to train the GAN. The inputs for this TorchGAN include the real dataset (flowers in our example or any other dataset which we have shown can easily be curated from online sources), dimension of the generated images, and the loss function. Here, we vary the input dimension for the TorchGAN from 16×16 to 128×128 to generate different patch sizes to feed to the classifier.

Attack Configuration and Success Measure. The adversary attempts to discover a naturalistic perturbation that can fool the classifier. We consider two different types of attacks: i) *targeted attack* where attackers aim to misclassify to a specific targeted label y_{target} ; and ii) *untargeted attack* where attackers only want to degrade the performance of the system, as in a denial-of-the-service attack, by fooling the classifier to recognize the object as $y \neq y_{\text{source}}$. In this work, we focus mainly on *targeted* attacks as it is more challenging to misclassify to a targeted label than simply cause a misclassification, and we chose the targeted class randomly.

All of these attacks are *universal* meaning that the attacker only needs one *single* patch to hijack the decision of the network to misclassify *any* input. One of the benefits of implementing a *universal* attack is that the attack is *input-agnostic* making it a strong attack regardless of the input, just as a backdoor in a conventional Trojan attack.

We used the Attack Success Rate (ASR) metric measured as the number of examples to successfully fool the target network over the total number of evaluated examples to evaluate the attack effectiveness.

Implementation Details. In our experiments, we use Pytorch [47] library for implementation and verify our method on a NVIDIA RTX2080 GPU. Since the main focus in this paper is on visual classification tasks, we assume that pixel values of inputs x are in the integer range of $[0, 255]$ or scaled to float range of $[0, 1]$ which correspond to the current practice of image processing and deep learning library. We used $\alpha = 0.01$, i.e. we changed the value of each latent value by 0.01 on each step. We selected the number of iterations to be 20. The

number of iterations and α values are chosen heuristically; sufficient for the adversarial example to reach the point of fooling the classifier while keeping the computational cost of experiments manageable.

VI. EVALUATION OF TNT EFFECTIVENESS

First, we intensively investigate the effectiveness of TnTs on ImageNet because of the fact that ImageNet classification benchmark has led to a great number of advances in image classification that some call “superhuman” [29]. We summarize our evaluations of TnT attacks on different scenarios:

- **Attack effectiveness on the entire ImageNet validation set.** Given the massive size of the dataset, previous works [34], [38] evaluated attack success on a sample of 100 random images (ImageNet-100). In addition, we want to evaluate the effectiveness of the discovered TnTs from a sample of 100 images the entire 50,000 images in the validation set (ImageNet-50K). Notably, *to the best of our knowledge, we are the first to evaluate the effectiveness of an adversarial patch on the entire validation set of ImageNet* (see Section VI-A).
- **Robustness to changes in patch locations.** Other attack methods such as LaVAN [34] have shown that an adversarial patch ASR could degrade significantly by *shifting the patch slightly in the image*. Therefore, we evaluate the robustness of the patch to location changes (see Section VI-B).
- **Black-box attack (Transferability of TnTs).** We assess the success of a black-box attacker. Hence, we evaluate the transferability of the known TnT on unknown networks trained with ImageNet (recall a black-box attacker has no knowledge about the attacked network) (see Section VI-C).
- **Attack effectiveness and generalization across other visual tasks.** Our attack method is generic; to demonstrate, we evaluate the generalization of the method on different visual classification tasks and datasets such as CIFAR10, GTSRB and PubFig (see Section VI-D).
- **Studies on the effect of random color and flower patches.** Since the TnTs occlude a part of the image, we want to understand if the phenomenon we observe can be explained by occlusion or a network biased to flowers or colors. Therefore, we evaluate the effect of occlusion by a patch as well as a random flower on the ImageNet classifier (see Section VI-E).





The TnTs we use in these experiments cover 10% to 20% of the input image, comparable with the patch size in AdvPatch [6] We investigate different patch sizes in Section VIII-B.

A. Attack Effectiveness on Entire ImageNet Validation Set

Since ImageNet is a huge dataset, deploying the algorithm on this dataset is time and power-consuming. Thus, initially, similar to previous works [34], [38] we only use a small number of samples (100 images) to find the TnTs. With the small sample size of 100 images, we successfully realize multiple different TnTs that fool the classifier with an attack

success rate of higher than 90% (on 100 randomly investigated images), while still maintaining the naturalism of the flower patch produced by the TnT Generator (Algorithm 1). Interestingly, we found *various* TnTs during our investigation, and in Table I we illustrate four examples (more examples and targeted labels are shown in later Sections). Each of the realized TnT has its own distinct features, but they all maintain the natural-looking of a flower; and once applied on *any* input image will misclassify the image to the targeted class $y = y_{\text{target}}$ through different pre-trained classifiers (*VGG-16*, *Inception-V3* or *WideResnet50*) of Pytorch.

TABLE I: Different TnTs found using ImageNet-100 for different models and their corresponding ASRs when applied to larger test sets.

Target	custard apple	arti-choke	pine apple	conch
Example				
ImageNet-100 (Attacker's test set)	96%	94%	96%	97%
Generalization across a large corpus of unseen data				
ImageNet-1000	93.6%	93.6%	95.1%	94.6%
ImageNet-50K	94.14%	94.51%	94.21%	95.13%
Network	<i>VGG-16</i>	<i>Inception-V3</i>	<i>WideResNet50</i>	

We also further verify the generalization of the discovered TnTs *found from a small sample of 100 images to attack a much larger sample size* in ImageNet. The results are shown in Table I. Surprisingly, the TnTs that we found in the 100-sample set generalize remarkably well to a bigger validation set (50K samples from the ImageNet validation set). For example, the TnT realized in *WideResNet50* with the ASR of 96% to fool *any* input image to the target *pineapple* class still maintains a high attack success rate of 95.1% for another 1000 random images (ImageNet-1000). To further verify the effectiveness of the attack, we deploy the TnTs found (from the 100-sample set) on the whole validation set of 50K images of ImageNet, and notably, we can still achieve a very high ASR of 94.21% on that whole validation set (50K images). Although there is a slight drop in the ASR of around 1.79% (from 96% to 94.21%), the ASR is still notably high.

Remark 3: This is a serious threat as attackers only need a small sample set to exploit the vulnerability; the attack is low cost to deploy, and the high attack success rate of TnTs found in a small sample set generalize well to unseen data outside of the attacker's test set.

B. Robustness to Changes in Patch Locations

Initially, we choose the trigger at the *lower-right corner*, however, as shown in other attacks with noisy patches, the location of the patch can dramatically reduce the attack success rate [34]. Based on this, we evaluate the robustness of TnTs to changes in its location. Given that a TnT is a naturalistic

TABLE II: Altering the location of a TnT increases ASR as it covers the main features of inputs. Notably, the TnT here realized from a small sample of 100 images (ImageNet-100), generalizes extremely well to larger populations to fool **any** inputs to the **targeted** class *conch*. An illustration of TnT locations are shown in Fig. 8 in Appendix B and results are from the *WideResNet50* pre-trained model from Pytorch [47].

Trigger Location	ImageNet-100 ASR	ImageNet-1000 ASR	ImageNet-50K ASR
lower right	92%	94.6%	95.13%
upper right	96%	96.6%	96.52%
upper left	99%	96.1%	95.61%
lower left	95%	94.9%	93.9%
center	92%	93.5%	94.56%
top	97%	96.5%	95.24%
bottom	96%	92.1%	93.31%
left	96%	93.0%	91.76%
right	97%	94.4%	93.56%

patch, in contrast to noisy pixels, shifting our patch to different locations increases the attack success rate as the trigger can now occlude potentially salient features of the benign inputs. Table II illustrates the effect of changing the location of the selected trigger on an input. By shifting the TnT to nine different locations (8 along borders and 1 at the center) on the attacker test set (ImageNet-100), the ASR increases from 92% (the lower-right corner position that we chose) to 96% (at the upper-right corner), and significantly increased to 99% (at the upper-left of images) since the patch possibly occluded the main feature of most original inputs.

Interestingly, these ASRs still hold strongly when we *assess generalizability by using TnTs discovered TnT from the small test set—ImageNet-100—to larger testing sets* of 1000 samples and 50K samples as detailed in Table II. This shows that our described ASR in the following sections (where we fix the patches at the lower-right corner on 100 random images) might not be the optimal ASR.

In addition, the consistently high ASR demonstrates our TnTs are not the result of occluding salient features of images; otherwise, we will see a variation of ASRs (high when occluding and low otherwise). Furthermore, this is the stronger and more difficult attack, a *targeted* attack; hence, occluding the main feature will not help fool the classifier to predict the targeted label. More importantly, *we also demonstrate the location invariance of TnTs in physical world deployments in real-time video demonstrations* (see Section IX).

Remark 4: TnTs are robust to changes in location.

C. Blackbox Attack—Transferability of TnTs

The attacks we have investigated thus far are under the *white-box attack* model; the attacker needs to have full knowledge of the target model under attack. With the high generalization of the TnTs shown in prior experiments, in this section, we aim to evaluate if TnTs discovered on a task can

be transferred to another network to mount an attack in a **black-box** setting.

A black-box attacker only needs to access a *Source* model to mount a white-box attack to extract TnTs. Then, the attacker can apply the discovered TnTs *blindly* to any other network that implements the same task and dataset. In this setting, we employ the Visual Recognition task on ImageNet implemented on a *Source* to attack a *Target* network. Notably, there are no qualitative results for black-box attacks in [34], [6] and the transferability for a targeted attack has been shown to be challenging in [4], following the setting in [41], we evaluate the transferability of our TnTs in an untargeted attack setting. The detailed results of the black-box attack are shown in Table XIII (**Appendix E**). The TnT realized in one model is highly transferable to another network. In general, we *observe* that TnT realized on a network less accurate for the task such as *SqueezeNet* (even with the source ASR of 99%) does not generalize well to other networks such as *WideResNet50* (with only 36% ASR). In contrast, TnTs realized from networks more accurate for the task such as *WideResNet50* (ASR of 97%) is highly generalizable and achieves high ASRs on other networks with ASRs of more than 60%.

Remark 5: We demonstrate multiple successful black-box attacks and confirm the TnTs discovered on a task in one network can be transferred to another network.

D. Attack Effectiveness and Generalization to Other Tasks

We further investigate the effectiveness of our TnTs on the following three contrasting tasks.

Scene classification (CIFAR10). The objective of the task is to generate the TnT flower that can fool the classifier to misdetect any scene with the TnT flower to be recognized as the target class, here we choose the random label (`car`) as the target class. The result shows that the generated TnT can misclassify *any* input to the targeted label with high a ASR up to 90.12% for the targeted attack.

Traffic sign recognition (GTSRB). This is a challenging task since the training dataset includes various *physical and environmental variations including different distances, lighting or occluding conditions*. Nevertheless, the discovered TnTs still achieve significantly ASRs of up to 95.75% in an untargeted setting, a significant increase in the ASR compared to 20.73% caused by the occlusion of same-size random color patches in Table III. A sample of TnTs realized and the corresponding ASRs are displayed in Table VII in Appendix B.

Face recognition (PubFig). In this task, most of the salient features learnt by a network are on a face and a network learns to ignore background information. *To fool the network to recognize as a specific target without occluding the main features of the face is both an interesting and challenging task*. Some of the results for untargeted attacks are shown in Table VII. For *targeted* attacks to a designated target such as Barack Obama, we implement TnTs covering 20% of the images at the lower-right corner (similar coverage to that used in our ablation studies in Section VI-E). We successfully fool the network with an ASR of up to 97.28% to misclassify

anyone with the TnT to a prominent targeted person (e.g. Barack Obama in our evaluation). Illustrations of successful TnTs are shown in Figure 3.

Remark 6: The vulnerability to TnT is generic and exists across multiple tasks.

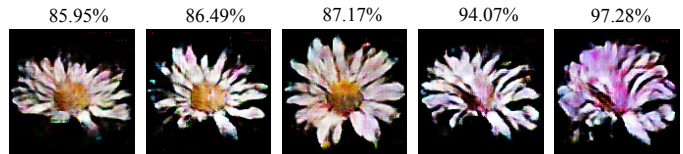


Fig. 3: TnTs realized in PubFig dataset and their corresponding ASRs for the *Face Recognition* task to impersonate *anyone* to the targeted person Barack Obama.

We employ the illustrated TnT in real-world attacks in a face recognition task in Fig. 6 of Section IX. Further results are shown in the **Appendix B**.

E. Can Occlusion or Network Bias Explain Attack Success?

In this section, we will examine the misclassification of the DNN system by using a random patch (random flowers or colors) to study the occlusion effect that a patch might have on the ASR of the visual task. In addition, we also utilize random flower patches to investigate if the behavior we observe can be explained by a *bias* in the network to flower images and to ascertain the possibility of randomly discovering a natural-looking flower that can fool the classifiers with a high ASR.

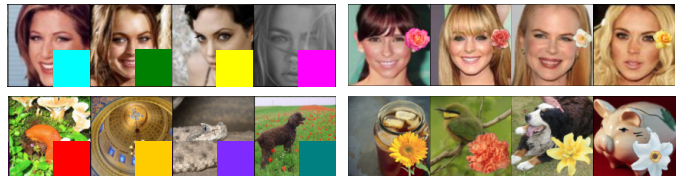


Fig. 4: Selective examples of random color and flower patches in our study on PubFig and ImageNet. The misclassification results caused by these patches are described in Table III.

TABLE III: Study results from affixing patches of either random colors or flowers to the test samples of each dataset. As observed, the success rate for an attacker employing such *simple tricks is fairly low*.

	Normal mis-classification	Random color patches	Random flower patches
CIFAR10	9.46%	25.99% \pm 0.432%	21.66% \pm 0.553%
GTSRB	3.23%	20.73% \pm 0.564%	18.71% \pm 0.665%
PubFig	5.26%	12.87% \pm 0.476%	6.02% \pm 0.031%
ImageNet-100	22%	22.72% \pm 0.109%	24.62% \pm 0.305%
ImageNet-1000	21.1%	24.85% \pm 0.15%	30.49% \pm 0.378%

Patch Size. *To examine the potential effects, all of the color and flower patches we selected will occlude the largest possible region we intentionally selected for our attack method (around 20% of the images).*

Random Colour Patches. We use 256 random color patches and digitally stamp the patch on the input (with the method described in the Equation 1) as a trigger at the lower-right corner of the image with the purpose of not occluding the main object of the image but assess the misclassification caused by the color patch.

Random Flowers Patches. We follow the approach with color patches but use randomly selected flowers drawn from our flower data set used for training the Generator. Particularly, we use 256 randomly chosen flowers and measure the attack success rate that a flower patch can cause on the classifier.

Results. In Table III, we reported the mean and standard deviation of the attack success rate for the patches for different tasks. Overall, the ASR achieved is significantly lower than the attack success rates demonstrated with TnTs (see Tables I and VII). Importantly, we observe a low standard deviation across all tasks; indicating that there is no special color or flower patch that can achieve a significantly high ASR compared to others. However, these results are far from a desirably high ASR to become a real threat, however, our investigation demonstrates that it is challenging to exploit a natural-looking patch while fooling the network with high ASRs. Notably, this low ASR is for an easy *untargeted* misclassification; hence, the ASR for the *targeted* attack is even much lower.

Remark 7: We demonstrate that exploiting a natural-looking patch to fool a network is a challenging task and the phenomenon we observed cannot be explained by occlusion or a network biased by flowers or colors; consequently, our attack method is an effective approach to realize such TnTs.

VII. EVALUATION OF TnT NATURALISM

In this section, we investigate the naturalism of the generated TnTs. We acknowledge that measuring naturalism is a challenging task, and there is no solid metric and definition fit for the purpose. However, following the studies in [64], [62], [5], [31], we consider measuring human perception of naturalism through user studies². We adopt the *Naturalistic Score* measure and the procedure in [31] to evaluate human perception of naturalism through two user studies (Study 1 and Study 2). For a robust evaluation, compared to previous studies [31] employing 10s of users, we conducted a *large cohort user study with 250 participants for each study*.






In user Study 1, we used a set of 9 patch images; i) 3 patches generated by LaVAN [34]; ii) 3 generated by AdvPatch [6]; and iii) 3 TnTs. All the patches are placed in random order and shown to participants (see Appendix D). The participants were asked to vote on each patch that looks natural to them. Then, we calculate the naturalistic score of each patch based on the percentage of participants’ votes. The aim of Study 1 is to measure the naturalism of TnTs compared to patches in previous attacks. The results in Table IV demonstrate our TnTs to have significantly higher naturalistic scores compared with the baselines.

²We followed Human Research Ethics Committee approval process, the study is considered ‘negligible risk’ and is exempt from ethical review.

In the second user study (Study 2), we randomly placed 3 of our generated TnTs together with 3 real flower images collected from Google Images [32] and asked participants to vote for the images that looks natural to them. The aim of this study is to measure the absolute naturalistic score of TnTs when compared with actual flower images. The results in Study 2 in Table IV demonstrate that our TnTs, synthetic images, can often be comparable to real flower images.

These results demonstrate our TnTs are more naturalistic and look significantly less malicious compared with prior works. Whilst our study was focused on evaluating the naturalism of the patch, we acknowledge that flowers may not look “in-place” in all settings and scenes. Automating the process of blending patches into a scene for an attack is another challenging problem, we will discuss further in Section XII. Notably, our method, whilst demonstrated on flowers, is generic and allows an adversary to self select natural objects beyond flowers as triggers, and thus facilitates the construction of natural-looking patches with the ability to blend better into the attacker-chosen setting.

TABLE IV: User studies to evaluate the Naturalistic Score of TnTs in comparison to other baseline patches. The Naturalistic Score is the percentage of participants’ votes. Complete details of patches used in the user studies are in Appendix D.

	Study 1 (250 participants)			Study 2 (250 participants)	
	LaVAN	AdvPatch	Ours	Google Images	Ours
Illustrative examples					
Naturalistic Scores (%)	0.4	12.0	90.0	97.2	89.6

VIII. GENERALIZATION TO ADVERSARIAL PATCH ATTACKS AND COMPARISON WITH PRIOR ATTACKS

To expand the scope of the attack for cases where there is no human involvement in the decision loop—for fully autonomous systems—and where stealth in a physical deployment is not an objective, we consider removing the naturalism constraint on the attack method. Therefore, we take a further step to let the TnT Generator G_θ learn the adversarial features from the classifier; thus, demonstrating the generic nature of our attack method—*i.e.* once we remove the naturalism constraint, our attack can generate conventional adversarial patches.

Attack Methodology. The proposed attack is detailed in *Algorithm 2* in Appendix C. We call this alternation an Adversarial Patch Generator since the Generator, after updating, learns the mapping from the latent vector \mathbf{z} to generate adversarial patches. More specifically, instead of searching in the latent space to find the TnT as illustrated in Fig. 2, now, we allow the Generator G_θ to be updated and learnt from the gradient back-propagation from the loss $\ell(\mathbf{x}', y_{\text{target}}, y_{\text{source}})$ to become the Adversarial Patch Generator ($G_{\theta'}$). This allows the Generator to learn the adversarial features and generate multiple adversarial triggers with different mappings from \mathbf{z} .

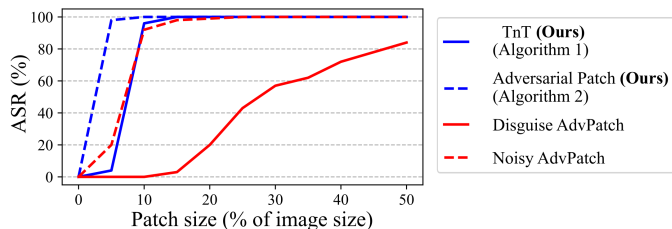


Fig. 5: Investigating Attack Success Rate (ASR) and patch size. Our TnT is comparable with the Noisy AdvPatch [6], and significantly better than the Disguise AdvPatch [6]. *Our Adversarial Patch ASR outperforms both Noisy and Disguise counterparts.*

A. Comparing to LaVAN: Smallest (Noisy) Adversarial Patch

To show the effectiveness of our patch attack, we opted to compare with LaVAN since the study achieved the *smallest* state-of-the-art patch results. We evaluate our patches with the same 14 targeted labels reported in LaVAN [34] on 100 random images. As expected, the patches of only 2% of the input image size from our Adversarial Patch Generator easily achieve 100% ASR on 100 randomly sampled images from ImageNet. The results in Table V demonstrate that our patches achieve a much higher ASR across the 14 targeted label on both settings of low and high confidence scores with a large margin of up to 44%.

TABLE V: Comparing ASR with LaVAN [34]—smallest state-of-the-art (noisy) adversarial patches of 2% of input image size—using the 14 targeted labels used in LaVAN on the *Inception-V3* network. Our patches achieve higher ASR in both settings; high and low confidence scores. Detailed examples of labels are shown in Appendix F.

	LaVAN	Ours
Average results across the 14 targeted labels from 100 ImageNet images		
ASR (conf \geq 0.9)	28.3%	72.9%
ASR (conf $<$ 0.9)	74.1%	98.1%

Remark 8: An attacker can trade-off naturalism to achieve significantly higher attack success rates compared with state-of-the-art adversarial patch attacks. Our results validates the generality of our attack method.

B. Comparing to AdvPatch: A Method to Disguise the Appearance of a Target Class in a Patch

Here, we compare with Disguise AdvPatch and Noisy AdvPatch from [6] on the same *VGG-16* network trained for the ImageNet task. We selected this study because: i) Disguise AdvPatch represents the efforts from [6] to hide the target class visible in the patch by disguising the noise patterns in another object while the Noisy AdvPatch is a noise-pattern patch visibly revealing the target class—*toaster*; ii) the method allows the generation of patches of different sizes. As shown in the Figure 5: i) efficacy of our TnT—the naturalistic patch—is comparable with the *Noisy AdvPatch*; and ii) significantly more effective than the Disguise AdvPatch

aiming to hide the true target class revealed to a human observer in the Noisy AdvPatch.

Notably, our naturalistic TnT patch is highly effective when the patch size is larger than 10% of the input image. To maintain the high ASR with smaller patches, we need to sacrifice some of the naturalism (Algorithm 2). Then, we can observe nearly 100% ASR in Figure 5 with a patch size of nearly 5% of the input image; now, the ASR is significantly higher than the noisy counterpart, Noisy AdvPatch.

IX. PHYSICAL WORLD DEPLOYMENTS

An advantage of a TnT is the ability for an adversary to easily print and deploy the attack in a scene in the physical world to fool a deep perception system. In this experiment, we print our TnTs and deploy *targeted* physical attacks for the ImageNet and PubFig tasks.

Attack Settings. Following the practices in the physical adversarial attack work in [36], we saved our triggers (TnTs) and patches (from the Adversarial Patch Generator) as .PNG files and printed the triggers and patches using a printer with a resolution of 300 *dpi* to maintain the pixel quality. Examples of the triggers and patches are shown in Figure 6 and 7. We validate the effectiveness of the attack in various physical world settings with complex backgrounds, as displayed in Figure. 6, by using a commercial camera of a smartphone to capture the scene—we used an iPhone 6S. We produced videos from our experiments to illustrate the effectiveness of the TnTs and patches in the physical world (<https://TnTattacks.github.io/>). In the videos, we also experiment with the robustness of our physical TnT in different *locations, scaling, lighting conditions, camera angles and positions, and so on.*

Results. Our results demonstrate that TnTs are robust to harsh physical-world conditions with *more than 90% of the images in frames successfully fooling the network and being recognized as the targeted class.* More detailed experiments of physical attacks are in videos accessed through the website.

We hypothesize the robustness of our TnT in the physical context is due to the fact that the patches are derived from a natural image distribution and are *universal or input-agnostic.* This allows the patch to potentially become invariant to various difficult conditions and suitable for deployment in physical world scenarios. Notably, we only apply the simplest method of physical deployment—printing without any modification to offset the printer quality and loss of pixels as in [51]; hence, the ASR and naturalism can potentially be improved by applying more robust adversarial printing techniques [51].

Remark 9: TnTs survive the harsh conditions in the physical world to pose a practical and realistic threat.

X. ATTACK EFFECTIVENESS AGAINST PATCH DEFENSES

The rise in adversarial patch attacks has led to the emergence of defenses—both empirical and provable methods [61], [18], [28], [13], [12], [37], [14], [17], [60], [49]. In this section, we evaluate the effectiveness of our TnTs against

both certified and empirical defenses against patch attacks. We defer details of the experimental setup and evaluation to Appendix A and summarize the key results here.

Against Empirically Robust Networks. In Table VI, we summarize the results for current state-of-the-art empirical defenses against physically realizable patch attacks. We report on DOA in [60], and the defense focusing on location optimization adversarial patches [49]. For the defended DOA network [60], when attacking the network with TnTs of the same patch size used in training (11×11), the robustness dropped from 80.43% to 13.78%. For the strongest defense, AT-FullLO network in [49], the robustness dropped from 72.2% to 11.02% under our TnT attack. In addition, because DOA was reported to be effective for sizes as large as 20% of the input [60], we increased the patch size of our TnTs to 20%. Although not reported in the table, the robustness of the network reduced significantly to 5.6%. Our hypothesis is that the TnTs, constructed under the naturalistic constraint, are forced to exploit vulnerabilities of the network to changes in natural features. These changes are different from those in adversarial patch exploration where the constraint is only on the spatial region occupied by the patch. Hence, a defense method built on adversarial training with adversarial patches as in [60], [49] is found to be largely ineffective against TnTs focusing on a different attack vector. This finding is also akin to the observations reported by Wu et al. [60] where the authors found that traditional defense methods such as adversarial training and randomized smoothing against unbounded spatial perturbation of adversarial examples are not robust against adversarial patches focusing on a *different attack vector*.

Although, deployed only digitally, we also conducted targeted LaVAN patch attacks based on the objective provided in the LaVAN study [34] to attack both defense methods for comparison. Albeit possessing the same capability and goal as our TnTs, the results in Table VI show that, unlike TnTs, LaVAN patches are ineffective against the empirically robust defense methods. Notably, the robustness to targeted attacks



Fig. 7: Physical deployment of TnTs generated from the TnT Generator targeting different classes (shown on top in red) in the ImageNet task.



Fig. 6: Various settings employed for the physical world attacks to impersonate ‘Barack Obama’. Results demonstrate our TnT is effective, even under different, complex, physical-world settings ranging from indoor to outdoor with different lighting conditions. The network recognizes the person with the TnT to be Barack Obama with high confidence.

from LaVAN is higher than that from untargeted Adversarial Patch attacks. The higher result, we hypothesize, is because the objective of misguiding the robust networks to a desired targeted label is more challenging. In summary, these results demonstrate: i) the effectiveness of TnTs, even against state-of-the-art patch defenses; and ii) that TnT attacks are an emerging new threat against DNNs.

Against Provably Robust Networks. For completeness and to assess the state-of-practice of current provable defenses, we evaluate the robustness of certified defenses after acknowledging the threat from TnTs—building the defense methods for TnT attacks. Given that the current state-of-the-art provable defenses are only effective against smaller patch sizes, the certified defenses achieved very low provable robustness; only up to 3.52% of inputs from 50,000 validation images can be certified in the best case (see Table VI). The reason is because, certifying against a larger patch, such as our TnTs or AdvPatch [6], requires certifying that a correct prediction can be made for a specific input, potentially tainted by a larger adversarial patch, in the presence of the defense method. In PatchGuard [6], this requires operating under larger masked regions in the feature space. Consequently, the provable defense must make predictions from the aggregation of the remaining features (not masked); leading to lower performance as well as certified robustness. Hence, provable defenses can only certify a smaller numbers of test inputs, i.e. achieve lower provable robustness, for larger adversarial patches. Therefore, an adversary can circumvent these defenses with a larger patch, such as our TnTs or even AdvPatch in [6]. Consequently, we observe TnTs attacks to still pose a realistic threat, even against *the strongest* defenses.

TABLE VI: Effectiveness of TnT attacks against robust defense methods (\uparrow *Robustness* is better for defenses).

Networks	Clean Acc	Provable Robustness ¹		
BagNet [21]	49.56%	0%		
Mask-BagNet [61]	49.65%	0.85%		
DS [37]	44.36%	3.52%		
Mask-DS [61]	39.77%	3.01%		
<i>Empirical Robustness</i>				
		Adversarial Patch ^{2,3}	LaVAN ³	TnT ³
DOA [60]	86.5%	80.43%	84.64%	13.78%
AT-FullLO [49]	90.44%	72.2%	78.44%	11.02%

¹Provable Robustness is attack-agnostic.

²Patches used are generated based on the code provided in studies [60] and [49]. These provide a baseline for comparisons with TnTs.

³The Adversarial Patch, LaVAN and TnT patches are of the same size: 11×11 .

XI. RELATED WORK

We describe prior work that focuses on *universal* perturbations and adversarial patch attacks and other *physically de-*

ployed adversarial attacks; further, we also compare our attack method with other GAN-based adversarial attack methods.

UAP and Adversarial Patch Attacks. Moosavi-Dezfooli et al. [41] showed the existence of a *universal* adversarial perturbation (UAP) in DNNs for image classification tasks, which is a *unique noise tensor* that when added to any input fools the classifier to mount an untargeted attack. The authors have also shown that there are multiple UAPs in a DNN, which can be transferable to another network architecture or for the same task. In order to deploy a *universal* adversarial attack in a real-world setting, Brown et al. [6] developed a spatially bounded adversarial patch (AdvPatch) to place in the scenes of ImageNet samples to fool the classifier to recognize objects as the `toaster` target class regardless of source inputs or locations. Karmon et al. extend this attack further to search for localized (or bounded) and visible adversarial noise in LaVAN [34] but the aim is to look for blind spots in a Deep Neural Network instead of physically realizable patches.

Other Physically Deployed Adversarial Attacks. Kukarin et al. [36] demonstrated that input-specific adversarial examples can also be deployed in the physical world in an untargeted attack if printed out and carefully cropped. Recently, Evtimov et al. [24] showed that specially crafted perturbations constrained to sticker shapes can fool a Traffic Sign recognition task once stuck to a Stop sign; while Athalye et al. [2] carefully crafted adversarial perturbations constrained to 3D objects to fool a DNN in the physical world. Different from ours, these adversarial examples are designed to work on a specific input (a specific Stop sign or 3D object), while our method is highly generalizable and the TnT can be printed out and attached to *any* input to work in the physical world (see Section IX). Concurrent studies [31], [56] have also attempted to realize naturalistic adversarial patch attacks, but against object detectors with the adversarial objective of causing a *misclassification of human objects* in a scene. To the best of our knowledge, our study remains the first to investigate universal, naturalistic, adversarial patches across a variety of classification tasks to misclassify *any* input to the attacker-desired *target label*. Further, TnTs are shown to be transferable and generalizable to tasks, models, and adversarial patch attacks.

GAN-based Adversarial Example Attacks. Researchers have investigated GAN-based structures to generate adversarial examples, such as [4], [15], [33], [66], [7]. Particularly, the authors in [66] train a GAN and an additional Inverter network to generate *full-size, fake images* that are able to flip the predicted label or mount an *untargeted attack*. Notably, these studies resemble the investigation of an adversarial example objective—input-dependent or noisy perturbation-based distortions added to an input *covering* the whole image to mount an *untargeted attack*. Different from this line of work, we rely on a *pre-trained* generator, and search the latent space \mathbf{z} to discover a type of spatially bounded adversarial example; a *patch* that is *physically realizable and universal (input-agnostic)*. These attributes ease the process of deployment in the physical world to mount *targeted* attacks.

GAN-based Adversarial Patch Attacks. Sharif et al. [52] apply a GAN-based method to generate spatially bounded phys-

ical adversarial sunglasses to impersonate a targeted person in a face recognition task (PubFig). The GAN-based method employs iterative training with feedback from the classifier under attack, which results in generating *noisy perturbations* constrained to the sunglasses mask. Notably, the sunglasses can be expected to occlude the main features of a face. In contrast, we do not alter the Generator for the TnT attacks (see Alg. 1) and are able to generate naturalistic patches by traversing through the latent space of the generator. Notably, the resulting patches can be successfully placed away from salient features of the input image.

PS-GAN [38] proposes to utilize a GAN-based structure to find patches with naturalism. The major differences with our work are: *i)* the attack is *untargeted* compared to both *targeted* and *untargeted* attacks capable with ours; *ii)* PS-GAN is *input-dependent* hence, an attacker needs to mount the attack in different ways for different inputs, which is harder to deploy in real-world scenarios; we address this problem using a *universal* or *input-agnostic* patch where *any* input will be misclassified when a TnT is applied; *iii)* PS-GAN patch is placed in the main context of the image, which can occlude main features; *iv)* method applies image-to-image translation using an encoder-decoder generator to translate an *existing* natural patch to an adversarial counterpart while ours learns the natural image distribution and approach a naturalistic adversarial patch. Similar to universal adversarial patches LaVAN [34] and AdvPatch [6]), we seek to be location invariant and attack method can lead to less malicious-looking and more powerful (higher ASR) attacks.

XII. DISCUSSION AND CONCLUSION

Are TnTs a formidable threat? We have validated through extensive experiments that natural-looking patches can successfully be used to fool Deep Neural Networks with high attack success rates. We have shown that TnTs are effective against multiple state-of-the-art classifiers ranging from widely used *WideResNet50* in the Large-Scale Visual Recognition task of ImageNet dataset to VGG-face models in the face recognition task of *PubFig* dataset in both *targeted* and *untargeted* attacks. TnTs can possess: *i)* the naturalism achievable with with triggers used in Trojan attack methods; and *ii)* the generalization and transferability of *adversarial examples* to other networks. This raises safety and security concerns regarding already deployed DNNs as well as future DNN deployments where attackers can use inconspicuous natural-looking object patches to misguide neural network systems without tampering with the model and risking discovery.

Are we limited to using flower triggers? In order to generate TnTs in our paper, we need at least one distribution of natural objects. In our experiments, we utilized flowers as explained. However, our attack generation method can generalize to any object selected by an attacker. Importantly, as we have demonstrated, it is both easy and low cost to obtain unlabelled datasets to generate TnTs.

What are potential avenues for mitigating the threat? We believe our work opens a new venue for further research into understanding the vulnerabilities of DNN systems. We believe our attack can be used for the *good* by providing

a method for not only discovering vulnerabilities of DNN models but to generate sample inputs to improve the robustness and trustworthiness of DNN models. Notably, inspired from other related work [60], [49], incorporating the emerging threat of TnTs within adversarial training may be a potential avenue. However, as we eluded to above, TnTs are not limited to flowers—only an illustrative naturalistic subset from other possible objects. Thus, incorporating flower TnTs within adversarial training might be effective against flower TnTs, but a more carefully-crafted natural set different from flower TnTs might still defeat such a defence. Thus, a general defence against TnTs appear to be an open problem.

Future Work. The TnTs generated whilst maintaining naturalism requires the patch to be approximately from 10% of the input image. Trading-off naturalistic features allows the trigger size to be 2% of the input image. We leave the task of reducing the size of the trigger whilst maintaining naturalistic features of TnTs for future research.

In order to generate TnTs, we need at least one distribution of natural objects. In our experiments, we utilized flowers as explained. However, natural objects are not limited to flowers and can be any object selected by an attacker. Importantly, it is both easy and low cost to obtain unlabelled datasets to generate TnTs (we used open-source, freely downloadable images to curate the flower dataset). We leave the investigation and exploration of other natural objects for future work.

It is an interesting research question to consider blending our patches into the scene to make the attack even further inconspicuous and stealthy. However, it is a challenging task since scenes captured evolves dynamically. Using naturalistic adversarial patches like ours is a first step to approach this problem, i.e. natural-looking flower patches are more likely to appear in a scene and are less malicious looking compared to noisy perturbation based patches in a scene. Methods to blend the universal adversarial patches into the scene require further research and we leave it for future work. Notably, our method, whilst demonstrated on flowers, allows an adversary to self select natural objects beyond flowers as triggers, and thus facilitates the construction of natural-looking patches with the ability to blend into the attacker-chosen settings.

Furthermore, our discussions here lead to the conclusion that learning a robust method against TnTs is non-trivial and challenging; and lead to an open problem worthwhile exploring in future research.

REFERENCES

- [1] Syed Muhammad Anwar et al. Medical image analysis using convolutional neural networks: a review. *Journal of Medical Systems*, 2018.
- [2] Anish Athalye, Logan Engstrom, Andrew Ilyas, and Kevin Kwok. Synthesizing robust adversarial examples. In *ICML*, 2018.
- [3] Eugene Bagdasaryan and Vitaly Shmatikov. Blind backdoors in deep learning models. In *USENIX Security*, 2021.
- [4] Shumeet Baluja and Ian Fischer. Learning to attack: Adversarial transformation networks. In *AAAI*, 2018.
- [5] Anand Bhattat et al. Unrestricted adversarial examples via semantic manipulation. In *ICLR*, 2020.
- [6] Tom B. Brown, Dandelion Mané, Aurko Roy, Martín Abadi, and Justin Gilmer. Adversarial patch. In *NIPS Workshop*, 2017.
- [7] Nicholas Carlini and Hany Farid. Evading deepfake-image detectors with white-and black-box attacks. In *CVPR Workshops*, 2020.
- [8] Nicholas Carlini, Matthew Jagielski, and Ilya Mironov. Cryptanalytic extraction of neural network models. In *CRYPTO*, 2020.
- [9] Nicholas Carlini and David Wagner. Adversarial examples are not easily detected: Bypassing ten detection methods. In *ACM Workshop on Artificial Intelligence and Security*, 2017.
- [10] Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *IEEE S&P*, 2017.
- [11] C. Chen, A. Seff, A. Kornhauser, and J. Xiao. DeepDriving: Learning affordance for direct perception in autonomous driving. In *ICCV*, 2015.
- [12] Ping-Yeh Chiang et al. Certified defenses for adversarial patches. In *ICLR*, 2020.
- [13] Edward Chou et al. Sentinet: Detecting localized universal attacks against deep learning systems. In *IEEE S&P Workshops*, 2020.
- [14] Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. Certified adversarial robustness via randomized smoothing. In *ICML*, 2019.
- [15] Chaowei Xiao et al. Generating adversarial examples with adversarial networks. In *IJCAI*, 2018.
- [16] J. Stallkamp et al. Man vs. computer: Benchmarking machine learning algorithms for traffic sign recognition. *Neural Networks*, 2012.
- [17] Matthew Mirman et al. Differentiable abstract interpretation for provably robust neural networks. In *ICML*, 2018.
- [18] Muzammal Naseer et al. Local gradients smoothing: Defense against localized adversarial attacks. In *WACV*, 2019.
- [19] O. Russakovsky et al. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*, 2015.
- [20] Tianyu Gu et al. Badnets: Identifying vulnerabilities in the machine learning model supply chain. 1708.06733, 2017.
- [21] Wieland Brendel et al. Approximating CNNs with bag-of-local-features models works surprisingly well on ImageNet. In *ICLR*, 2019.
- [22] Xinyun Chen et al. Targeted backdoor attacks on deep learning systems using data poisoning, 2017.
- [23] Yingqi Liu et al. Trojaning attack on neural networks. In *NDSS*, 2018.
- [24] Kevin Eykholt et al. Robust physical-world attacks on deep learning visual classification. In *CVPR*, 2018.
- [25] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *ICLR*, 2015.
- [26] Ian Goodfellow et al. Generative adversarial nets. In *NIPS*, 2014.
- [27] Ishaan Gulrajani et al. Improved training of wasserstein gans. In *NIPS*, 2017.
- [28] Jamie Hayes. On visible adversarial perturbations & digital watermarking. In *CVPR Workshops*, 2018.
- [29] Kaiping He et al. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *ICCV*, 2015.
- [30] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. In *CVPR*, 2021.
- [31] Yu-Chih-Tuan Hu et al. Naturalistic physical adversarial patch for object detectors. In *CVPR*, 2021.
- [32] Google Inc. *Google Images*, 2020 (accessed July 3, 2020).
- [33] Sargan Jandial et al. AdvGAN++: Harnessing latent layers for adversary generation. In *ICCV Workshops*, 2019.
- [34] Danny Karmon, Daniel Zoran, and Yoav Goldberg. LaVAN: Localized and visible adversarial noise. In *ICML*, 2018.
- [35] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-10 (canadian institute for advanced research).
- [36] Alexey Kurakin, Ian Goodfellow, and Samy Bengio. Adversarial examples in the physical world. *ICLR Workshop*, 2017.
- [37] Alexander Levine and Soheil Feizi. (De)randomized smoothing for certifiable defense against patch attacks. In *NeurIPS*, 2020.
- [38] Aishan Liu et al. Perceptual-sensitive gan for generating adversarial patches. In *AAAI*, 2019.
- [39] Jiajun Lu et al. No need to worry about adversarial examples in object detection in autonomous vehicles. In *CVPR Workshops*, 2017.
- [40] Aleksander Madry et al. Towards deep learning models resistant to adversarial attacks. In *ICLR*, 2018.
- [41] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. Universal adversarial perturbations. In *CVPR*, 2017.
- [42] OpenCV. *Image Thresholding*, 2020 (accessed June 25, 2020). https://opencv-python-tutroals.readthedocs.io/en/latest/py_tutorials/py_imgproc/py_thresholding/py_thresholding.html.
- [43] Avik Pal and Aniket Das. TorchGAN: A Flexible Framework for GAN Training and Evaluation, 2019.
- [44] Nicolas Papernot et al. Transferability in machine learning: from phenomena to black-box attacks using adversarial samples. *arXiv:1605.07277*, 2016.
- [45] Nicolas Papernot et al. Practical black-box attacks against machine learning. In *AsiaCCS*, 2017.

- [46] O. M. Parkhi, A. Vedaldi, and A. Zisserman. Deep face recognition. In *British Machine Vision Conference*, 2015.
- [47] A. Paszke et al. Pytorch: An imperative style, high-performance deep learning library. In *NeurIPS*, 2019.
- [48] Nicolas Pinto et al. Scaling up biologically-inspired computer vision: A case study in unconstrained face recognition on facebook. 2011.
- [49] Sukrut Rao, David Stutz, and Bernt Schiele. Adversarial training against location-optimized adversarial patches. 2020.
- [50] David Rolnick and Konrad Kording. Reverse-engineering deep ReLU networks. In *ICML*, 2020.
- [51] Mahmood Sharif et al. Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition. In *CCS*, 2016.
- [52] Mahmood Sharif et al. A general framework for adversarial examples with objectives. *ACM Transactions on Privacy and Security*, 2019.
- [53] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.
- [54] Christian Szegedy et al. Intriguing properties of neural networks. *arXiv:1312.6199*, 2013.
- [55] Yaniv Taigman et al. DeepFace: Closing the gap to human-level performance in face verification. In *CVPR*, 2014.
- [56] Jia Tan, Nan Ji, Haidong Xie, and Xueshuang Xiang. Legitimate adversarial patches: Evading human eyes and detection models in the physical world. In *Proceedings of the 29th ACM International Conference on Multimedia*, 2021.
- [57] Florian Tramèr, Fan Zhang, Ari Juels, Michael K Reiter, and Thomas Ristenpart. Stealing machine learning models via prediction apis. In *USENIX Security*, 2016.
- [58] Florian Tramèr et al. On adaptive attacks to adversarial example defenses. *arXiv preprint arXiv:2002.08347*, 2020.
- [59] Bolun Wang et al. Neural cleanse: Identifying and mitigating backdoor attacks in neural networks. In *IEEE S&P*, 2019.
- [60] Tong Wu, Liang Tong, and Yevgeniy Vorobeychik. Defending against physically realizable attacks on image classification. In *ICLR*, 2020.
- [61] Chong Xiang et al. Patchguard: A provably robust defense against adversarial patches via small receptive fields and masking. In *USENIX Security*, 2021.
- [62] Chaowei Xiao et al. Spatially transformed adversarial examples. In *ICLR*, 2018.
- [63] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In *British Machine Vision Conference*, 2016.
- [64] Richard Zhang et al. Colorful image colorization. In *ECCV*, 2016.
- [65] Xinyang Zhang et al. Interpretable deep learning under fire. In *USENIX Security*, 2020.
- [66] Zhengli Zhao, Dheeru Dua, and Sameer Singh. Generating natural adversarial examples. In *ICLR*, 2018.

APPENDIX A

ATTACK EFFECTIVENESS AGAINST PATCH DEFENSES

A. Against Empirically Robust Networks

In this section, we evaluate our TnTs against empirically robust networks. As shown in [65], [9], [58], empirical defenses such as [13], [18], [28] are usually vulnerable to adaptive attackers once they are aware of the working mechanisms of the defenses. Notably, recently, Wu et al. [60] highlighted that the conventional methods to improve the robustness against adversarial examples such as adversarial training and randomized smoothing showed limited effectiveness against physically realizable adversarial attacks, and proposed an approach named Defense against Occlusion Attacks (DOA) to defend against these physically realizable adversarial patch attacks. This approach is also developed in Rao et al. [49] but the authors jointly optimize patch values and location. These two state-of-the-art empirical defenses [60], [49] are the most relevant defense method against our attacks. Hence, we employ these defenses in this section and assess the effectiveness of our TnTs against these robust defenses.

Experiment Setup. We use the pre-trained robust Resnet110 networks provided by Wu et al. [60] on Github³ for the CIFAR-10 dataset. This network was trained on a patch size of 11×11 . Notably, as reported in the paper, even a network trained with a smaller patch size can defend against a patch size as large as 20% of the images (*i.e.* cover all of our attacks). For Rao et al. [49], we train the given robust ResNet-20 model from scratch following the default parameters of the strongest proposed defense (AT-FullLO)⁴ on CIFAR-10 dataset. To make a fair comparison with [60], we also train with the same mask size of 11×11 . To assess the effectiveness of our attack against these robust networks, we deploy our TnTs to achieve the challenging task of fooling the network to misclassify *any* images to the targeted label (*car*). First, we evaluate the robustness of the network against our TnTs for the same patch size used in training (11×11). However, to stretch out the robustness of the defended network, we also evaluate the network against a larger patch size of TnT around 20% since it was reported to be effective for sizes as large as 20% of the input [60]. In addition, we employed the adversarial patches in [60], [49] to compare with our TnT attacks as shown in Table VI. LaVAN patches employed are of the same size, 11×11 , and aim to hijack the robust network decision to the target class—*car*—similar to our TnTs.

Metrics. We report the *Clean Acc*—the Accuracy on benign inputs, and *Empirical Robustness*—the percentage of input images with patches that are correctly classified or the performance of the network under an attack.

B. Against Provably Robust Networks

A provable defense is the strongest defense and could potentially block and eliminate the adversarial patches completely. In this section, we evaluate the robustness of our attack

against the multiple state-of-the-art (SOTA) provable defenses including BagNet [21], Derandomized Smoothing [37], and recently the improved versions of those defenses named PatchGuard [61] demonstrating superior performance compared with other provable methods [12], [37], [14], [17]. Notably, the PatchGuard method relies on the small receptive fields and robust masking to eliminate the adversarial effects of malicious patches and generate provable robustness for the defended system.

Experiment Setup. In this experiment, we use the same networks and configurations as in [61] of BagNet [21] with receptive fields of 17×17 and de-randomized smoothed ResNet (DS) [37] evaluated on ImageNet. Then, we also apply Robust Masking defense in [61] to generate Mask-BN and Mask-DS provable robust networks, respectively. The provable Robustness results in Table VI is evaluated on the entire validation set (50,000 images) of ImageNet using a mask size of 10% of the input image size.

Metrics. We use two metrics in this experiment: i) *Clean Acc*—the Accuracy obtained on benign inputs from the testing set; and ii) *Provable Robustness*—the percentage of the images in the clean testing set that are able to certified (*i.e.* no attack is expected to succeed). We report the results in Table VI.

APPENDIX B

ATTACK EFFECTIVENESS AND GENERALIZATION TO OTHER TASKS

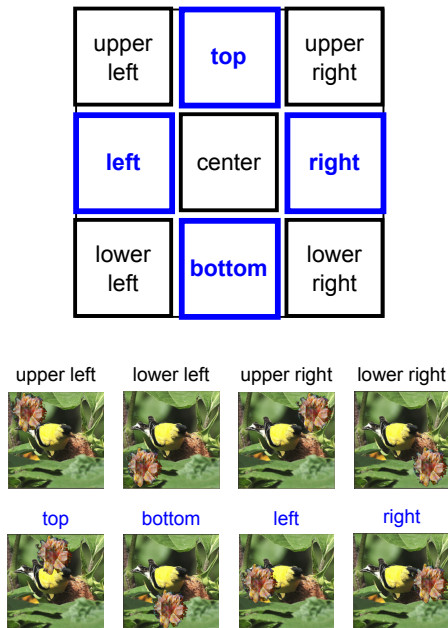


Fig. 8: An illustration of trigger locations around the border.

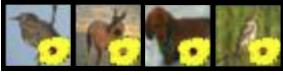

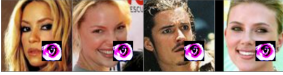
Effectiveness of Patch Locations. From Table II, we can see that with 8 locations around the border of the input images (excluding the center location), TnTs achieved the maximum ASR of 96.52% (at upper-right corner) and the minimum ASR of 91.76% (at left corner). Throughout all of 8 border locations, TnTs achieved a high mean of 94.4% with a low

³<https://github.com/tongwu2020/phattacks>

⁴<https://github.com/sukrutrao/Adversarial-Patch-Training>

variation of only 1.53% showing the effectiveness of our TnTs even at border locations.

TABLE VII: Illustrative examples of TnTs found in different visual classification tasks and their corresponding ASRs

Dataset	ASR	Examples	Target
CIFAR10	90.12%		car
GTSRB	95.75%		untargeted
PubFig	95.14%		untargeted

APPENDIX C

DETAILED INFORMATION ON DATASETS, MODEL ARCHITECTURES AND TRAINING CONFIGURATIONS

TABLE VIII: Networks used for the classification tasks

Task/Dataset	# of Labels	# of Training Images	# of Testing Images	Model Architecture
CIFAR10[35]	10	50,000	10,000	6 Conv + 2 Dense
GTSRB[16]	43	39,288	12,630	7 Conv + 2 Dense
PubFig[48]	170	48,498	12,322	13 Conv + 3 Dense (VGG-16)
ImageNet[19]	1,000	⁵	50,000	WideResNet50 [63]

We describe the datasets used in our studies below.

- **Large Scale Visual Recognition** (ImageNet [19]). ImageNet is a highly popular real-world dataset with a million high-resolution images of a large variety of objects used for training state-of-the-art deep perception models. The goal is to recognize visual scenes among 1,000 different classes. This is one of the most popular dataset in computer vision for benchmarks state-of-the-art models. In this task, we utilized state-of-the-art *pre-trained* models available from Pytorch Deep Learning library [47]; notably, these models are used as base models by many Machine Learning practitioners for transfer learning to build systems for different visual tasks.

Additionally, to demonstrate the generalization of our method, we also evaluate on 3 other visual classification tasks: i) Scene Classification (CIFAR10); ii) Traffic Sign Recognition (GTSRB); and iii) Face Recognition (PubFig).

- **Scene Classification** (CIFAR10 [35]). This is a widely used task and dataset with images of size 32×32 and we used a similar network to that implemented in the IEEE S&P [59] study.
- **Traffic Sign Classification** (GTSRB [16]). German Traffic Sign Benchmark dataset is commonly used to evaluate vulnerabilities of DNNs as it is related to autonomous

driving and safety concerns. The goal is to recognize 43 different traffic signs of size 32×32 to simulate a scenario in self-driving cars. The network we used follows the VGG [53] network structure.

- **Face Recognition** (PubFig [48]). Public Figures Face dataset is a large, real-world dataset with high-resolution images of large variations in pose, lighting, and expression. The goal is to recognize the faces of public figures. In this task, we leverage transfer learning from a publicly available pre-trained model based on a complex 16-layer VGG-Face model from the work of [46] and fine-tune the last 6 layers.

TABLE IX: Model Architecture for VGGFace2

Layer Type	# of Channels	Filter Size	Stride	Activation
Conv	64	3	1	ReLU
Conv	64	3	1	ReLU
MaxPool	64	2	2	-
Conv	128	3	1	ReLU
Conv	128	3	1	ReLU
MaxPool	128	2	2	-
Conv	256	3	1	ReLU
Conv	256	3	1	ReLU
Conv	256	3	1	ReLU
MaxPool	256	2	2	-
Conv	512	3	1	ReLU
Conv	512	3	1	ReLU
Conv	512	3	1	ReLU
MaxPool	512	2	2	-
Conv	512	3	1	ReLU
Conv	512	3	1	ReLU
Conv	512	3	1	ReLU
MaxPool	512	2	2	-
FC	4096	-	-	ReLU
FC	4096	-	-	ReLU
FC	170	-	-	Softmax

TABLE X: Model Architecture for GTSRB

Layer Type	# of Channels	Filter Size	Stride	Activation
Conv	128	3	1	ReLU
Conv	128	3	1	ReLU
MaxPool	128	2	2	-
Conv	256	3	1	ReLU
Conv	256	3	1	ReLU
MaxPool	256	2	2	-
Conv	512	3	1	ReLU
Conv	512	3	1	ReLU
MaxPool	512	2	2	-
Conv	1024	3	1	ReLU
MaxPool	1024	2	2	-
FC	1024	-	-	ReLU
FC	10	-	-	Softmax

TABLE XI: Model Architecture for CIFAR-10. FC: fully connected layer.

Layer Type	# of Channels	Filter Size	Stride	Activation
Conv	128	3	1	ReLU
Conv	128	3	1	ReLU
MaxPool	128	2	2	-
Conv	256	3	1	ReLU
Conv	256	3	1	ReLU
MaxPool	256	2	2	-
Conv	512	3	1	ReLU
Conv	512	3	1	ReLU
MaxPool	512	2	2	-
FC	1024	-	-	ReLU
FC	10	-	-	Softmax

⁵In our work, we utilized pre-trained models on the ImageNet dataset from Pytorch [47].

TABLE XII: Dataset and Training Configuration

Task/Dataset	# of Labels	Input Size	Training Set Size	Testing Set Size	Training Configurations
CIFAR-10	10	$32 \times 32 \times 3$	50,000	10,000	epochs=100, batch=32, optimizer=Adam, lr=0.001
GTSRB	43	$32 \times 32 \times 3$	35,288	12,630	epochs=25, batch=32, optimizer=Adam, lr=0.001
PubFig	83	$224 \times 224 \times 3$	11,070	2,768	epochs=30, batch=32, optimizer=Adam, lr=0.001. The last 4 layers are fine-tuned during training.
ImageNet	1,000	$224 \times 224 \times 3$	-	50,000	We utilize pre-trained models available from Pytorch [47] for the task.

Algorithm 2: Adversarial Patch Generator Process

1 **Inputs:** a batch of images $\{\mathbf{x}^{(i)}\}_{i=1}^m$ with batch size m , source label $\{y_{\text{source}}^{(i)}\}_{i=1}^m$, target label $\{y_{\text{target}}^{(i)}\}_{i=1}^m$, model p_M , latent vector \mathbf{z} , the hyper-parameter λ to balance the loss, natural generator G_θ , and the thresholds to detect TnT τ_{batch} , τ_{val} for batch and validation set respectively;

2 **Initialization;**

3 **while** $ASR < \tau_{\text{val}}$ **do**

4 Sample a batch of images $\{\mathbf{x}^{(i)}\}_{i=1}^m$;

5 Sample a latent variable $\mathbf{z} \sim p(\mathbf{z})$;

6 $\delta = G_\theta(\mathbf{z})$ ▷ a flower patch;

7 Generate the mask \mathbf{m} based on δ ;

8 $\delta' \leftarrow \text{bgremoval}(\delta, \mathbf{m})$ ▷ Background removal;

9 **for** $i = 1, \dots, m$ **do**

10 $\mathbf{x}'^{(i)} = (1 - \mathbf{m}) \odot \mathbf{x}^{(i)} + \mathbf{m} \odot \delta'$;

11 $y_{\text{argmax}}^{(i)} = \arg \max_y p_M(y | \mathbf{x}'^{(i)})$;

12 **if** $y_{\text{argmax}}^{(i)} = y_{\text{target}}$ **then**

13 $f_{\text{fool}} = f_{\text{fool}} + 1$;

14 **end if**

15 **end for**

16 $L = \ell(\{\mathbf{x}'^{(i)}\}_{i=1}^m, \{y_{\text{target}}^{(i)}\}_{i=1}^m) - \lambda \ell(\{\mathbf{x}'^{(i)}\}_{i=1}^m, \{y_{\text{source}}^{(i)}\}_{i=1}^m)$;

17 $\theta \leftarrow \text{Adam}(\nabla_{\theta} L, \theta, \alpha, \beta_1, \beta_2)$;

18 **if** $f_{\text{fool}} > \tau_{\text{batch}}$ **then**

19 Sample a latent variable $\mathbf{z} \sim p(\mathbf{z})$;

20 $\delta = G_\theta(\mathbf{z})$;

21 Test this δ for the whole validation set \mathcal{X}_{val} to get ASR ;

22 **if** $ASR \geq \tau_{\text{val}}$ **then**

23 Complete update Generator, save the latest state as Adversarial Patch Generator $G_{\theta'}$;

24 **end if**

25 **end if**

26 **end while**

APPENDIX D

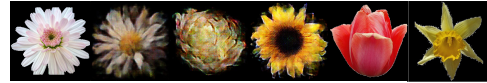
DETAILED INFORMATION ON USER STUDIES

We adopted the Naturalistic Score and follow the procedure described in [31] to evaluate the human perception of our TnTs' naturalism using two study designs. As in [31], naturalism was defined as the *Naturalistic Score* of each patch calculated based on the percentage of participants' votes for the patch. In contrast to [31], employing 24 participants, we conduct a *large* cohort user study with 250 participants for each of the following studies:

- **User study 1.** The aim of this study is to measure and compare the naturalistic score of our TnTs with previous attacks. We used a set of 9 patch images; i) 3 patches generated by LAVAN [34]; ii) 3 generated by AdvPatch [6]; and iii) 3 TnTs. All the patches are placed in random order and shown to participants and the participants were asked to vote on each patch that looks natural to them (see Figure 9a).
- **User study 2.** The aim of this study is to measure the absolute naturalistic score of our TnTs when compared with real flower images. We randomly placed 3 of our generated TnTs together with 3 real flower images collected from Google Images [32] and asked participants to vote on the images that looks natural to them (see Figure 9b).



(a) Study 1









(b) Study 2

Fig. 9: An instance of the random ordering of patch images used in the two user studies, 250 participants participated in each study.

APPENDIX E
DETAILED INFORMATION ON BLACK-BOX ATTACK

Notably, to the best of our knowledge, we are the first to report qualitatively the untargeted black-box attack success rates where the patches—TnT in our attack—do not occlude the main salient features of the images (Table XIII).

TABLE XIII: Black-box attack, the transferability of TnT from a model to other models on ImageNet dataset in an untargeted setting. ASRs are observed on 100 random images (network performance on the task is given in parenthesis).

Example		Target					
		WideResNet50	Inception V3	ResNet18	SqueezeNet 1.0	VGG 16	MnasNet
	WideResNet50 (Acc: 78.51%)	97%	77%	67%	77%	78%	63%
	Inception V3 (Acc: 77.45%)	46%	91%	51%	80%	66%	57%
	ResNet18 (Acc: 69.76%)	45%	47%	80%	64%	59%	54%
	SqueezeNet 1.0 (Acc: 58.1%)	36%	42%	51%	99%	47%	48%
	VGG 16 (Acc: 71.59%)	38%	49%	49%	70%	91%	49%
	MnasNet (Acc: 73.51%)	47%	63%	59%	74%	56%	87%

APPENDIX F
ILLUSTRATIVE EXAMPLES OF ADVERSARIAL PATCH GENERATOR OUTPUTS FOR TARGET LABELS IN LAVAN

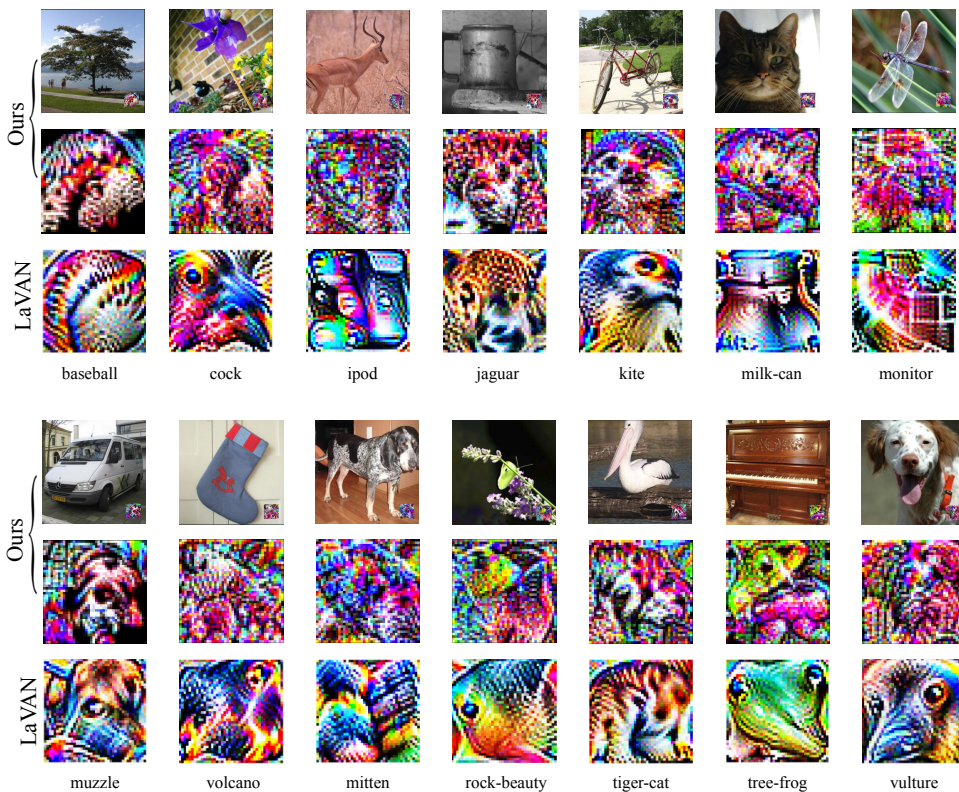


Fig. 10: Generated patches from our Adversarial Patch Generator for the 14 different targeted labels in LaVAN. **The 1st row:** Our adversarial patches in the scene. **The 2nd row:** Our adversarial patches were rescaled to image size for display purposes. **The 3rd row:** adversarial patches from LaVAN for the same targeted label.