# Hardly Perceptible Trojan Attack against Neural Networks with Bit Flips

Jiawang Bai<sup>1</sup>, Kuofeng Gao<sup>1</sup>, Dihong Gong<sup>2</sup>, Shu-Tao Xia<sup>1,3, $\boxtimes$ </sup>, Zhifeng Li<sup>2, $\boxtimes$ </sup>, and Wei Liu<sup>2, $\boxtimes$ </sup>

<sup>1</sup> Tsinghua Shenzhen International Graduate School, Tsinghua University <sup>2</sup> Data Platform, Tencent

<sup>3</sup> Research Center of Artificial Intelligence, Peng Cheng Laboratory {bjw19,gkf21}@mails.tsinghua.edu.cn; gongdihong@gmail.com; xiast@sz.tsinghua.edu.cn; michaelzfli@tencent.com; wl2223@columbia.edu

Abstract. The security of deep neural networks (DNNs) has attracted increasing attention due to their widespread use in various applications. Recently, the deployed DNNs have been demonstrated to be vulnerable to Trojan attacks, which manipulate model parameters with bit flips to inject a hidden behavior and activate it by a specific trigger pattern. However, all existing Trojan attacks adopt noticeable patch-based triggers (e.g., a square pattern), making them perceptible to humans and easy to be spotted by machines. In this paper, we present a novel attack, namely hardly perceptible Trojan attack (HPT). HPT crafts hardly perceptible Trojan images by utilizing the additive noise and per-pixel flow field to tweak the pixel values and positions of the original images, respectively. To achieve superior attack performance, we propose to jointly optimize bit flips, additive noise, and flow field. Since the weight bits of the DNNs are binary, this problem is very hard to be solved. We handle the binary constraint with equivalent replacement and provide an effective optimization algorithm. Extensive experiments on CIFAR-10, SVHN, and ImageNet datasets show that the proposed HPT can generate hardly perceptible Trojan images, while achieving comparable or better attack performance compared to the state-of-the-art methods. The code is available at: https://github.com/jiawangbai/HPT.

## 1 Introduction

Although deep neural networks (DNNs) have been showing state-of-the-art performances in various complex tasks, such as image classification [47,43,17], facial recognition [50,28,8,38], and object detection [13,42], prior studies have revealed their vulnerability against diverse attacks [15,34,46,16,36,27,44,48,5]. One such attack is the *Trojan attack* [32] happening in the deployment stage, in which an attacker manipulates a DNN to inject a hidden behavior called Trojan. The Trojan can only be activated by the specific trigger pattern.

 $<sup>\</sup>boxtimes$  Corresponding authors.



Fig. 1: Visualization of Trojan images from TrojanNN [32], TBT [40], and HPT. '*Trans.*' indicates the use of transparent trigger proposed in [32]. Note that ProFlip [6] uses the same trigger pattern (a square pattern) as TBT.

Trojan attacks on a deployed DNN alter the model parameters in the memory using bit flip techniques, e.g., Row Hammer Attack [20,51], but do not tamper with the training pipeline and have no extra forward or backward calculation during inference. Then, the attacked DNN makes a target prediction on the inputs with the trigger, while behaving normally on clean samples [32,40,6]. These dangerous properties pose severe threats to DNN-based applications after model deployment. Therefore, it is necessary to study the Trojan attacks on the deployed DNNs in order to recognize their flaws and solve related security risks.

A Trojan attack generally is composed of two subroutines: critical bits identification and specific trigger generation. Previous works [39,6] made efforts towards reducing the number of bit flips by developing search strategies. After flipping the identified bits, the attacker can activate the hidden behavior using the Trojan images, which are any images embedded with a specific patch-based trigger, such as a watermark pattern in [32] or a square pattern in [40,6]. However, due to these unnatural and noticeable triggers, the Trojan images can be easily spotted by humans [37,9] and machines. For example, we construct a simple linear classifier to distinguish clean images and Trojan images crafted by TBT [40] based on their Grad-CAM heatmaps [45] on ImageNet, resulting in a 98.0% success rate. The transparent trigger [32] was proposed to reduce the perceptibility, but leading to lower attack performance. Hence, how to inject less perceptible Trojan with superior attack performance is a challenging problem.

To address the aforementioned problems, we propose the hardly perceptible Trojan attack (HPT). Instead of applying the patch-based trigger predefined by a mask, HPT tweaks the pixel values and positions of the original images to craft Trojan images. Specifically, we modify pixel values by adding the pixel-wise noise inspired by adversarial examples [49,15,34] and change pixel positions by per-pixel flow field [10,19,63]. As shown in Figure 1, Trojan images of HPT are less perceptible and harder to be distinguished from original images. It will be further demonstrated in the human perceptual study in Section 4.3.

Since the value of each weight bit is '0' or '1', we cast each bit as a binary variable. Integrating Trojan images generation and critical bits identification, we formulate the proposed HPT as a mixed integer programming (MIP) problem, i.e., jointly optimizing bit flips, additive noise, and flow field. Moreover, we constrain the modification on image pixels and the number of bit flips to yield a hardly perceptible and practical attack [51,62]. To solve this MIP problem, we reformulate it as an equivalent continuous problem [56] and present an effective

Table 1: Summary of attributes of TrojanNN [32], TBT [40], ProFlip [6], and HPT. '*Trans.*' indicates the use of transparent trigger proposed in [32] and 'ASR' denotes the attack success rate.

| sees the account  | access rat   | 0.                       |                            |
|-------------------|--------------|--------------------------|----------------------------|
| Method            | High ASR     | A Small Set of Bit Flips | Hardly Perceptible Trigger |
| TrojanNN          | $\checkmark$ |                          |                            |
| TrojanNN (Trans.) |              |                          | $\checkmark$               |
| TBT               | $\checkmark$ |                          |                            |
| ProFlip           | $\checkmark$ | $\checkmark$             |                            |
| HPT (Ours)        | $\checkmark$ | $\checkmark$             | $\checkmark$               |

optimization algorithm based on the standard alternating direction method of multipliers (ADMM) [12,4]. We conduct extensive experiments on CIFAR-10, SVHN, and ImageNet with 8-bit quantized ResNet-18 and VGG-16 architectures following [40,6], which shows that HPT is hardly perceptible and achieves superior attack performance. The attributes of compared methods and HPT are summarized in Table 1.

The contributions are summarized as follows:

- For the first time, we improve Trojan attacks on the deployed DNNs to be strong and hardly perceptible. We investigate the use of modifying the pixel values and positions of the original images to craft Trojan images.
- We formulate the proposed HPT as a constrained MIP problem to jointly optimize bit flips, additive noise, and flow field, and further provide an effective optimization algorithm.
- Finally, HPT obtains both hardly perceptible Trojan images and promising attack performance, e.g., an attack success rate of 95% with only 10 bit flips out of 88 million bits in attacking ResNet-18 on ImageNet.

## 2 Related Works and Preliminaries

Attacks on the Deployed DNNs. Recently, since DNNs have been widely applied to security-critical tasks, e.g., facial recognition [14,55,53,59], their security in the deployment stage has received extensive attention. Previous works assume that the attacker can modify the weight parameters of the deployed DNNs to achieve some malicious purposes, e.g., misclassifying certain samples [31,62]. Some physical techniques (e.g., Row Hammer Attack [20,51] and Laser Beam Attack [1,7]) can precisely flip bits ( ${}^{\circ}O \rightarrow {}^{\circ}1'$  or  ${}^{\circ}1'\rightarrow {}^{\circ}O'$ ) in the memory without accessing them. These techniques allow an attacker to attack a deployed DNN by modifying its bit representation [60,52,41]. For instance, Rakin et al. [39] presented that an extremely small number of bit flips can crush a fully functional DNN to a random output generator. After that, Bai et al. [3] proposed to attack a specified sample into a target class via flipping limited bits. To mitigate the bit flip-based attacks, some defense mechanisms [18,25,30,24] have been explored.

As a line of research, *Trojan attacks* [40,6] insert a hidden behavior in the DNN using bit flip techniques, which can be activated by a designed trigger.

#### 4 Jiawang Bai, Kuofeng Gao, Dihong Gong et al.



Fig. 2: Pipeline of the proposed HPT, where the target class is '*Goldfish*'. (a) Trojan Injection: optimizing  $\hat{\theta}$ ,  $\delta$ , and f jointly. (b) Inference: activating the hidden behaviour using the Trojan image.

Specifically, an attacked DNN will make wrong predictions on the inputs with the presence of the trigger, while behaving normally on original samples. Previous works proposed to generate a specified pattern as the trigger (e.g., a square pattern) [32,40,6]. Besides, to reduce the number of bit flips, heuristic strategies are designed to identify the critical bits, e.g., neural gradient ranking in [40] or progressive search algorithm in [6]. The results in [6] show that only a few bit flips yield a high attack success rate, which further raises the security concerns. **Quantized DNNs.** In the deployment stage, model quantization has been widely adopted to reduce the storage requirements and accelerate the inference speed [29,21]. In this paper, we adopt a layer-wise Q-bit uniform quantizer, which is identical to the Tensor-RT solution [33]. Given binary weight parameters  $\boldsymbol{\theta} \in \{0, 1\}^{N \times Q}$  of a Q-bit quantized DNN g, each parameter is represented and stored as a signed integer in two's complement representation, i.e.,  $\boldsymbol{\theta}_i = [\boldsymbol{\theta}_{i,Q}; ...; \boldsymbol{\theta}_{i,1}] \in \{0, 1\}^Q$ . For the l-th layer with the step size  $\Delta_l$ , the binary representation  $\boldsymbol{\theta}_i$  can be converted into a real number, as follows:

$$\boldsymbol{W}_{i} = (-2^{Q-1} \cdot \boldsymbol{\theta}_{i,Q} + \sum_{j=1}^{Q-1} 2^{j-1} \cdot \boldsymbol{\theta}_{i,j}) \cdot \boldsymbol{\Delta}_{l},$$
(1)

where  $\boldsymbol{W} \in \mathbb{R}^N$  denotes the floating-point weight parameters of the DNN. For clarity, hereafter  $\boldsymbol{\theta}$  is reshaped from the tensor to the vector, i.e.,  $\boldsymbol{\theta} \in \{0, 1\}^{NQ}$ . **Threat Model.** We consider the threat model used by previous bit flip-based Trojan attack studies [40,6]. The attacker knows the location of the victim DNN in the memory to implement precisely bit flips. We also assume that the attacker has a small set of clean data and knows the architecture and parameters of the victim DNN. Note that our attack does not have access to the training process nor the training data. During inference, the attacker can activate the injected Trojan by applying the generated trigger on the test samples.

## 3 Methodology

In this section, we firstly describe the hardly perceptible trigger used by HPT. We then introduce the problem formulation and present an effective optimization algorithm based on ADMM. Figure 2 shows the entire pipeline of HPT.

### 3.1 Hardly Perceptible Trigger

Let  $\boldsymbol{x} \in \mathcal{X}$  denote a clean image, and  $\hat{\boldsymbol{x}}$  denote its Trojan image, i.e.,  $\boldsymbol{x}$  with a specific trigger pattern.  $\mathcal{X} = [0, 1]^{HW \times C}$  is the image space, where H, W, C are the height, width, and channel for an image, respectively. To modify the pixel values of the clean image, we apply additive noise  $\boldsymbol{\delta}$  to the image  $\boldsymbol{x}$ . Inspired by the adversarial examples [23], HPT requires  $\boldsymbol{\delta} \in S_n$ , where  $S_n = \{\boldsymbol{\delta} : ||\boldsymbol{\delta}||_{\infty} \leq \epsilon\}$  and  $\epsilon$  denotes the maximum noise strength, so that  $\hat{\boldsymbol{x}}$  is perceptually indistinguishable from  $\boldsymbol{x}$ .

To formulate changes of the pixel positions, we use  $\mathbf{f} \in \mathbb{R}^{HW \times 2}$  to represent per-pixel flow field, where  $\mathbf{f}^{(i)} = (\Delta u^{(i)}, \Delta v^{(i)})$  denotes the amount of displacement in each channel for each pixel  $\hat{\mathbf{x}}^{(i)}$  within the Trojan image at the position  $(\hat{u}^{(i)}, \hat{v}^{(i)})$ . Thus, the value of  $\hat{\mathbf{x}}^{(i)}$  is sampled from position  $(u^{(i)}, v^{(i)}) = (\hat{u}^{(i)} + \Delta u^{(i)}, \hat{v}^{(i)} + \Delta v^{(i)})$  within the original image. Since the sampled position is not necessarily an integer value, we use the differentiable bilinear interpolation [63] to generate the Trojan image considering four neighboring pixels around the location  $(u^{(i)}, v^{(i)})$ , denoted by  $\mathcal{N}(u^{(i)}, v^{(i)})$ . To preserve high perceptual quality of the Trojan images, we enforce the local smoothness of the flow field  $\mathbf{f}$  based on the total variation [57,61]:

$$\mathcal{F}(\boldsymbol{f}) = \sum_{p}^{all \ pixels} \sum_{q \in \mathcal{N}(p)} \sqrt{||\Delta u^{(p)} - \Delta u^{(q)}||_2^2 + ||\Delta v^{(p)} - \Delta v^{(q)}||_2^2}.$$
 (2)

We introduce a hyper-parameter  $\kappa$  to restrict  $\mathcal{F}(\mathbf{f})$ , i.e.,  $\mathbf{f} \in \mathcal{S}_f$  where  $\mathcal{S}_f = \{\mathbf{f} : \mathcal{F}(\mathbf{f}) \leq \kappa\}$ .

Based on the additive noise  $\delta$  and the flow field f, each pixel before the clipping operation can be calculated as:

$$\hat{\boldsymbol{x}}^{(i)} = \sum_{q \in \mathcal{N}(u^{(i)}, v^{(i)})} (\boldsymbol{x}^{(q)} + \boldsymbol{\delta}^{(q)}) (1 - |u^{(i)} - u^{(q)}|) (1 - |v^{(i)} - v^{(q)}|).$$
(3)

We can craft the Trojan image  $\hat{x}$  by calculating each pixel of  $\hat{x}$  according to Eqn. (3) and performing the [0,1] clipping to ensure that it is in the image space, which is denoted as:

$$\hat{\boldsymbol{x}} = \mathcal{T}(\boldsymbol{x}; \boldsymbol{\delta}, \boldsymbol{f}). \tag{4}$$

Note that  $\hat{x}$  is differentiable with respect to  $\delta$  and f, enabling us to optimize them by the gradient method. We obtain  $\delta$  and f after the Trojan injection stage and apply them on any image to craft its Trojan image during inference.

### 3.2 Problem Formulation

Suppose that the victim DNN is the well-trained Q-bit quantized classifier  $g : \mathcal{X} \to \mathcal{Y}$ , where  $\mathcal{Y} = \{1, ..., K\}$  is the output space and K is the number of classes.  $\boldsymbol{\theta} \in \{0, 1\}^{NQ}$  is the original binary weight parameters and  $\hat{\boldsymbol{\theta}} \in \{0, 1\}^{NQ}$  is the  $\mathbf{6}$ 

attacked parameters. As aforementioned, the attacker has a small set of clean data  $D = \{(\boldsymbol{x}_i, y_i)\}_{i=1}^{M}$ . To keep the attack stealthiness, we reduce the influence on original samples by minimizing the below loss:

$$\mathcal{L}_{cle}(\hat{\boldsymbol{\theta}}) = \sum_{i=1}^{M} \ell(g(\boldsymbol{x}_i; \hat{\boldsymbol{\theta}}), y_i),$$
(5)

where  $g_j(\boldsymbol{x}_i; \hat{\boldsymbol{\theta}})$  indicates the posterior probability of the input with respect to class j and  $\ell(\cdot, \cdot)$  denotes the cross-entropy loss. Recall that the malicious purpose of Trojan attack is to classify the Trojan images into the target class t. To this end, we formulate this objective as:

$$\mathcal{L}_{tro}(\boldsymbol{\delta}, \boldsymbol{f}, \hat{\boldsymbol{\theta}}) = \sum_{i=1}^{M} \ell(g(\mathcal{T}(\boldsymbol{x}_i; \boldsymbol{\delta}, \boldsymbol{f}); \hat{\boldsymbol{\theta}}), t).$$
(6)

Aligning with previous Trojan attacks [40,6], reducing the number of bit flips is necessary to make the attack efficient and practical. HPT achieves this goal by restricting the Hamming distance between  $\theta$  and  $\hat{\theta}$  less than b. Considering the constraint on  $\delta$  and f, we formulate the objective function of HPT as follows:

$$\{\boldsymbol{\delta}^{*}, \boldsymbol{f}^{*}, \hat{\boldsymbol{\theta}}^{*}\} = \arg\min_{\boldsymbol{\delta}, \boldsymbol{f}, \hat{\boldsymbol{\theta}}} \mathcal{L}_{cle}(\hat{\boldsymbol{\theta}}) + \gamma \mathcal{L}_{tro}(\boldsymbol{\delta}, \boldsymbol{f}, \hat{\boldsymbol{\theta}}),$$
  
s.t.  $\boldsymbol{\delta} \in \mathcal{S}_{n}, \quad \boldsymbol{f} \in \mathcal{S}_{f}, \quad \hat{\boldsymbol{\theta}} \in \{0, 1\}^{NQ}, \quad d_{H}(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}) \leq b,$  (7)

where  $\gamma$  is the hyper-parameter to balance the two terms and  $d(\theta, \hat{\theta})$  computes the Hamming distance between original and attacked parameters. Problem (7) with the continuous variables  $\delta$  and f and the binary variable  $\hat{\theta}$  is a mixed integer programming, which is generally very difficult to solve. Here, inspired by a recent advanced work in integer programming [56], we equivalently replace the binary constraint with the intersection of two continuous constraints, as presented in Proposition 1.

**Proposition 1.** [56] Let  $\mathbf{1}_{NQ}$  denote the NQ-dimensional vector filled with all 1s. The binary set  $\{0,1\}^{NQ}$  can be replaced by the intersection between  $\mathcal{S}_b$  and  $\mathcal{S}_p$ , as follows:

$$\hat{\boldsymbol{\theta}} \in \{0,1\}^{NQ} \Leftrightarrow \hat{\boldsymbol{\theta}} \in (\mathcal{S}_b \cap \mathcal{S}_p),\tag{8}$$

where  $S_b = [0,1]^{NQ}$  indicates the box constraint and  $S_p = \{\hat{\theta} : ||\hat{\theta} - \frac{1}{2}\mathbf{1}_{NQ}||_2^2 = \frac{NQ}{4}\}$ indicates the  $\ell_2$ -sphere constraint.

Based on Proposition 1, we can equivalently reformulate Problem (7) as a continuous problem. Besides, we can calculate the Hamming distance using  $||\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}||_2^2$  for the binary vectors. We obtain the following reformulation:

$$\begin{aligned} \{\boldsymbol{\delta}^{*}, \boldsymbol{f}^{*}, \hat{\boldsymbol{\theta}}^{*}\} &= \arg\min_{\boldsymbol{\delta}, \boldsymbol{f}, \hat{\boldsymbol{\theta}}} \mathcal{L}_{cle}(\hat{\boldsymbol{\theta}}) + \gamma \mathcal{L}_{tro}(\boldsymbol{\delta}, \boldsymbol{f}, \hat{\boldsymbol{\theta}}), \\ s.t. \quad \boldsymbol{\delta} \in \mathcal{S}_{n}, \quad \boldsymbol{f} \in \mathcal{S}_{f}, \hat{\boldsymbol{\theta}} = \boldsymbol{z}_{1}, \quad \hat{\boldsymbol{\theta}} = \boldsymbol{z}_{2}, \quad ||\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}||_{2}^{2} - b + z_{3} = 0, \end{aligned}$$
(9)  
$$\boldsymbol{z}_{1} \in \mathcal{S}_{b}, \quad \boldsymbol{z}_{2} \in \mathcal{S}_{p}, \quad z_{3} \in \mathbb{R}^{+}. \end{aligned}$$

In Eqn. (9), we use two additional variables  $z_1$  and  $z_2$  to separate the two continuous constraints in Proposition 1, and transform  $||\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}||_2^2 \leq b$  into  $\{||\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}||_2^2 - b + z_3 = 0; z_3 \in \mathbb{R}^+\}$ . Problem (7) can now be optimized by the standard ADMM method.

### 3.3 Optimization

Using a penalty factor  $\rho > 0$ , the augmented Lagrangian function of Eqn. (9) is

$$L(\boldsymbol{\delta}, \boldsymbol{f}, \hat{\boldsymbol{\theta}}, \boldsymbol{z}_{1}, \boldsymbol{z}_{2}, \boldsymbol{z}_{3}, \boldsymbol{\lambda}_{1}, \boldsymbol{\lambda}_{2}, \boldsymbol{\lambda}_{3}) = \mathcal{L}_{cle}(\hat{\boldsymbol{\theta}}) + \gamma \mathcal{L}_{tro}(\boldsymbol{\delta}, \boldsymbol{f}, \hat{\boldsymbol{\theta}}) + \boldsymbol{\lambda}_{1}^{\top}(\hat{\boldsymbol{\theta}} - \boldsymbol{z}_{1}) + \boldsymbol{\lambda}_{2}^{\top}(\hat{\boldsymbol{\theta}} - \boldsymbol{z}_{2}) + \boldsymbol{\lambda}_{3}^{\top}(||\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}||_{2}^{2} - b + \boldsymbol{z}_{3})$$

$$+ \frac{\rho}{2} \left[ ||\hat{\boldsymbol{\theta}} - \boldsymbol{z}_{1}||_{2}^{2} + ||\hat{\boldsymbol{\theta}} - \boldsymbol{z}_{2}||_{2}^{2} + (||\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}||_{2}^{2} - b + \boldsymbol{z}_{3})^{2} \right]$$

$$+ \mathbb{I}_{\mathcal{S}_{n}}(\boldsymbol{\delta}) + \mathbb{I}_{\mathcal{S}_{f}}(\boldsymbol{f}) + \mathbb{I}_{\mathcal{S}_{b}}(\boldsymbol{z}_{1}) + \mathbb{I}_{\mathcal{S}_{p}}(\boldsymbol{z}_{2}) + \mathbb{I}_{\mathbb{R}^{+}}(\boldsymbol{z}_{3}),$$

$$(10)$$

where  $\lambda_1 \in \mathbb{R}^{NQ}$ ,  $\lambda_2 \in \mathbb{R}^{NQ}$ , and  $\lambda_3 > 0$  are Lagrange multipliers for the three equality constraints.  $\mathbb{I}_S(a) : a \to \{0, +\infty\}$  denotes the indicator function:  $\mathbb{I}_S(a) = 0$  if a belongs to a set S; otherwise,  $\mathbb{I}_S(a) = +\infty$ .

We alternatively update all variables as shown in Algorithm 1. We start from initializing all optimizable variables and the iteration index k (Line 2-4), where the initialization of  $\boldsymbol{\delta}^{[0]}$  and  $\boldsymbol{f}^{[0]}$  is specified later. We first update  $\boldsymbol{z}_1^{[k+1]}, \boldsymbol{z}_2^{[k+1]}$ , and  $\boldsymbol{z}_3^{[k+1]}$  with Eqn. (11)-(12) (Line 8-10), which project  $\boldsymbol{z}_1^{[k+1]}, \boldsymbol{z}_2^{[k+1]}$ , and  $\boldsymbol{z}_3^{[k+1]}$  into their feasible sets:

$$\Pi_{\mathcal{S}_b}(\boldsymbol{e}_1) = \min(1, \max(0, \boldsymbol{e}_1)), \tag{11}$$

$$\Pi_{\mathcal{S}_p}(\boldsymbol{e}_1) = \frac{(2\boldsymbol{e}_1 - \mathbf{1}_{NQ})\sqrt{NQ}}{||\boldsymbol{e}_1||} + \frac{\mathbf{1}_{NQ}}{2},\tag{12}$$

$$\Pi_{\mathbb{R}^+}(e) = \max(0, e),\tag{13}$$

where  $e_1 \in \mathbb{R}^{NQ}$  and  $e \in \mathbb{R}$ . Next, we update  $\delta^{[k+1]}, f^{[k+1]}$ , and  $\hat{\theta}^{[k+1]}$  by gradient descent with learning rates  $\alpha_{\delta}, \alpha_{f}$ , and  $\alpha_{\hat{\theta}}$ , respectively (Line 12-14). The projection functions for  $\delta^{[k+1]}$  and  $f^{[k+1]}$  are defined as:

$$\Pi_{\mathcal{S}_n}(\boldsymbol{e}_2) = \min(-\epsilon, \max(\boldsymbol{e}_2, \epsilon)), \qquad (14)$$

$$\Pi_{\mathcal{S}_f}(\boldsymbol{e}_3) = \frac{\boldsymbol{e}_3}{\mathcal{F}(\boldsymbol{e}_3)},\tag{15}$$

where  $\boldsymbol{e}_2 \in \mathbb{R}^{HW \times C}$  and  $\boldsymbol{e}_3 \in \mathbb{R}^{HW \times 2}$ . We then update the Lagrange multipliers  $\boldsymbol{\lambda}_1^{[k+1]}, \boldsymbol{\lambda}_2^{[k+1]}$ , and  $\boldsymbol{\lambda}_3^{[k+1]}$  using gradient ascent (Line 16-18). When both  $||\hat{\boldsymbol{\theta}} - \boldsymbol{z}_1||_2^2$  and  $||\hat{\boldsymbol{\theta}} - \boldsymbol{z}_2||_2^2$  are smaller than a preset threshold or the maximum number of iterations is reached, the optimization halts (Line 19).

Jiawang Bai, Kuofeng Gao, Dihong Gong et al.

Algorithm 1: ADMM for solving Problem (7)

**Input** : Victim DNN model q with binary weight parameters  $\theta$ , target class t, a small set of clean data  $D = \{(\boldsymbol{x}_i, y_i)\}_{i=1}^M$ . Output:  $\delta^*, f^*, \theta^*$ . 1 # Initialization **2** Initialize  $\boldsymbol{\delta}^{[0]}$  and  $\boldsymbol{f}^{[0]}$ ; **3** Let  $\hat{\boldsymbol{\theta}}^{[0]} \leftarrow \boldsymbol{\theta}, \boldsymbol{z}_{1}^{[0]} \leftarrow \boldsymbol{\theta}, \boldsymbol{z}_{2}^{[0]} \leftarrow \boldsymbol{\theta}, \boldsymbol{z}_{3}^{[0]} \leftarrow 0, \boldsymbol{\lambda}_{1}^{[0]} \leftarrow \boldsymbol{0}_{NQ}, \boldsymbol{\lambda}_{2}^{[0]} \leftarrow \boldsymbol{0}_{NQ}, \boldsymbol{\lambda}_{3}^{[0]} \leftarrow 0;$ 4 Set k = 0; 5 repeat 6 |  $\ \ \, \#$  Update  $z_1^{[k+1]}, z_2^{[k+1]}, z_3^{[k+1]}$ 7 until Stopping criterion is satisfied;  $\begin{array}{l} \mathbf{s} \hspace{0.1cm} \boldsymbol{z}_{1}^{[k+1]} \leftarrow \Pi_{\mathcal{S}_{b}}(\hat{\boldsymbol{\theta}}^{[k]} + \boldsymbol{\lambda}_{1}^{[k]}/\rho); \\ \mathbf{g} \hspace{0.1cm} \boldsymbol{z}_{2}^{[k+1]} \leftarrow \Pi_{\mathcal{S}_{p}}(\hat{\boldsymbol{\theta}}^{[k]} + \boldsymbol{\lambda}_{2}^{[k]}/\rho); \end{array}$  $\begin{array}{l} \mathbf{10} \ z_3^{[k+1]} \leftarrow \Pi_{\mathbb{R}^+}(-||\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}^{[k]}||_2^2 + b - \lambda_3^{[k]}/\rho); \\ \mathbf{11} \ \mathbf{\#} \ \mathbf{Update} \ \boldsymbol{\delta}^{[k+1]}, \boldsymbol{f}^{[k+1]}, \hat{\boldsymbol{\theta}}^{[k+1]}. \end{array}$ 12  $\boldsymbol{\delta}^{[k+1]} \leftarrow \Pi_{\mathcal{S}_n} (\boldsymbol{\delta}^{[k]} - \alpha_{\boldsymbol{\delta}} \cdot \partial L / \partial \boldsymbol{\delta});$ 13  $\boldsymbol{f}^{[k+1]} \leftarrow \Pi_{\mathcal{S}_f} (\boldsymbol{f}^{[k]} - \alpha_f \cdot \partial L / \partial \boldsymbol{f});$ 14  $\hat{\boldsymbol{\theta}}^{[k+1]} \leftarrow \hat{\boldsymbol{\theta}}^{[k]} - \alpha_{\hat{\boldsymbol{\theta}}} \cdot \partial L / \partial \hat{\boldsymbol{\theta}};$ 
$$\begin{split} & \mathbf{14} \ \boldsymbol{\theta}^{(k-1)} \leftarrow \boldsymbol{\theta}^{(k-1)} - \boldsymbol{\lambda}_{\boldsymbol{\theta}}^{(k-1)} \cdot \boldsymbol{DL} \\ & \mathbf{15} \ \ \mathbf{\#} \ \mathbf{Update} \ \boldsymbol{\lambda}_{1}^{[k+1]}, \mathbf{\lambda}_{2}^{[k+1]}, \mathbf{\lambda}_{3}^{[k+1]} \\ & \mathbf{16} \ \boldsymbol{\lambda}_{1}^{[k+1]} \leftarrow \boldsymbol{\lambda}_{1}^{[k]} + \rho(\hat{\boldsymbol{\theta}}^{[k+1]} - \boldsymbol{z}_{1}^{[k+1]}); \\ & \mathbf{17} \ \boldsymbol{\lambda}_{2}^{[k+1]} \leftarrow \boldsymbol{\lambda}_{2}^{[k]} + \rho(\hat{\boldsymbol{\theta}}^{[k+1]} - \boldsymbol{z}_{2}^{[k+1]}); \\ & \mathbf{18} \ \boldsymbol{\lambda}_{3}^{[k+1]} \leftarrow \boldsymbol{\lambda}_{3}^{[k]} + \rho(||\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}^{[k+1]}||_{2}^{2} - b + \boldsymbol{z}_{3}^{[k+1]}); \end{split}$$
**19**  $k \leftarrow k + 1;$ 20  $\delta^* \leftarrow \delta^{[k]}, f^* \leftarrow f^{[k]}, \hat{\theta}^* \leftarrow \hat{\theta}^{[k]};$ 21 return  $\delta^*, f^*, \hat{\theta}^*$ .

**Implementation Details.** We implement the optimization process with the following techniques. We initialize  $\delta^{[0]}$  and  $f^{[0]}$  by minimizing the loss defined as Eqn. (6) before joint optimization to stabilize the practical convergence. In the step for  $\hat{\theta}$ , we only update the parameters of the last layer and fix the others. We also update  $\delta^{[k+1]}$ ,  $f^{[k+1]}$ , and  $\hat{\theta}^{[k+1]}$  with multi-steps gradients during each iteration. Besides, as suggested in [56,26,11], increasing  $\rho$  from a smaller value to an upper bound can avoid the early stopping.

## 4 Experiments

#### 4.1 Setup

**Datasets and Target Models.** Following [40,6], we adopt three datasets: CIFAR-10 [22], SVHN [35], and ImageNet [43]. The attacker has 128 clean images for CIFAR-10 and SVHN and 256 clean images for ImageNet, respectively. Note that all attacks are performed using these clean images and evaluated on the whole test set. Following [40,6], we evaluate attacks on two popular network

8

| Dataset  | Model                | Method            | <b>PA-TA</b> (%) | $\mathbf{ASR}\ (\%)$ | $N_{\rm flip}$ |
|----------|----------------------|-------------------|------------------|----------------------|----------------|
|          |                      | TrojanNN          | 87.6             | 93.9                 | 19215          |
|          | DeeNet 19            | TrojanNN (Trans.) | 75.5             | 73.5                 | 20160          |
|          | TA: 04.807           | TBT               | 87.5             | 90.2                 | 540            |
|          | 1A: 94.070           | TBT $(Trans.)$    | 71.4             | 56.6                 | 548            |
| CIEAR 10 |                      | HPT (Ours)        | 94.7             | 94.1                 | 12             |
| CIFAR-10 |                      | TrojanNN          | 85.5             | 82.5                 | 16400          |
|          | VCC 16               | TrojanNN (Trans.) | 69.6             | 59.8                 | 15386          |
|          | TA: 93.2%            | TBT               | 80.7             | 83.2                 | 601            |
|          |                      | TBT $(Trans.)$    | 70.5             | 53.9                 | 583            |
|          |                      | HPT (Ours)        | 93.1             | 91.1                 | 6              |
|          |                      | TrojanNN          | 76.0             | 82.5                 | 17330          |
|          | VCC 16               | TrojanNN (Trans.) | 59.9             | 71.7                 | 18355          |
| SVHN     | VGG-10<br>TA: 06.207 | TBT               | 67.9             | 60.1                 | 576            |
|          | IA: 90.370           | TBT $(Trans.)$    | 57.9             | 54.8                 | 546            |
|          |                      | HPT (Ours)        | 94.2             | 78.0                 | 26             |
|          |                      | TrojanNN          | 47.6             | 100.0                | 155550         |
|          | DecNet 19            | TrojanNN (Trans.) | 47.4             | 96.8                 | 304744         |
| ImageNet | TAL CO FOZ           | TBT               | 68.8             | 100.0                | 611            |
|          | 1A: 09.5%            | TBT $(Trans.)$    | 64.1             | 88.6                 | 594            |
|          |                      | HPT (Ours)        | 68.6             | 95.2                 | 10             |

Table 2: Performance comparison between TrojanNN [32], TBT [40], and HPT. 'Trans.' indicates the use of transparent trigger proposed in [32]. The target class t is set as 0.

architectures: ResNet-18 [17] and VGG-16 [47], with a quantization level of 8-bit (see Appendix B for results of HPT in attacking 4-bit quantized DNNs). In the below experiments, the target class t is set as 0 unless otherwise specified.

**Parameter Settings.** To balance the attack performance and human perception,  $\epsilon$  is set as 0.04 on all datasets, and  $\kappa$  is set as 0.01 on CIFAR-10 and SVHN, and 0.005 on ImageNet. We initialize  $\delta^{[0]}$  and  $f^{[0]}$  by minimizing loss defined as Eqn. (6) for 500 iterations with the learning rate 0.01 on CIFAR-10 and SVHN, and 1,000 iterations with the learning rate 0.1 on ImageNet. Other parameter settings can be found in Appendix A.

**Evaluation Metrics.** To measure the effect on clean images, we compare original test accuracy (TA) with post-attack test accuracy (PA-TA), defined as the accuracy of testing on clean images for the original and attacked DNN, respectively. Attack success rate (ASR) denotes the percentage of Trojan images samples classified to the target class by the attacked DNN.  $N_{flip}$  is the number of bit flips, i.e., the hamming distance between original and attacked parameters. A more successful attack can achieve a higher PA-TA and ASR, while less  $N_{flip}$ . **Compared Methods.** HPT is compared to TrojanNN [32], TBT [40], and ProFlip [6] in our experiments. We also apply the transparent trigger [32] on TrojanNN and TBT, denoted as '*Trans.*'. The watermark in [32] is chosen as the trigger shape for TrojanNN. We use the same trigger pattern in [32,40], i.e., a square pattern located at the bottom right of the image. The trigger size of all compared methods is measured by the proportion of input replaced by the trigger, which is configured as the default value used in [32,40], i.e., 9.76% on CIFAR-10 and SVHN, and 10.62% on ImageNet. We use the open-sourced code



Fig. 3: Visualization of clean and Trojan images generated by different methods. Top left: CIFAR-10; Bottom left: SVHN; Right: ImageNet. Note that ProFlip uses the same trigger pattern (a square pattern) as TBT. By comparison, the triggers of Trojan images generated by HPT are most natural and unnoticeable.



Fig. 4: Comparison of ProFlip [6] and HPT: (a) Attack performance. The target class t is set as 2; (b) Visualization of clean and Trojan images generated by ProFlip and HPT. The example of ProFlip is from [6].

of TBT and implement TrojanNN following [40] to make the comparison fair. We compare the results of ProFlip reported in [6].

### 4.2 Attack Results

We compare the attack performance of HPT with TrojanNN and TBT in Table 2. We can observe that TrojanNN performs a large number of bit flips in all cases, due to no limitation on the parameter modification. TBT reduces the number of bit flips to about 600 using the proposed bit search algorithm. Among them, HPT achieves the least number of bit flips with a higher or at least competitive PA-TA and ASR. It is worth noting that the transparent trigger leads to a lower PA-TA and ASR for both TrojanNN and TBT. For example, in attacking ResNet-18 on CIFAR-10, TBT only achieves a 56.6% ASR, compared to the 90.2% ASR without applying the transparent trigger, which indicates that it is difficult to perform Trojan attacks with the less perceptible trigger.

The results of ProFlip and HPT with the target class t = 2 are shown in Figure 4(a). The number of clean samples is set as 256 for all cases. ProFlip is the most state-of-the-art method based on the well-designed progressive search

Table 3: Scores of human perceptual Table 4: Performance of HPT with study (ranging from 1 to 5). A higher different target classes t in attacking score corresponds to less perceptible Tro- ResNet-18 on CIFAR-10. n images į٤

| an mag   | ges.     |          |     |          |     |  | t 1 | (07) | (07) | Nflip | t | (07) | (07) | Nflip |
|----------|----------|----------|-----|----------|-----|--|-----|------|------|-------|---|------|------|-------|
|          |          | TrojanNN | -   | TBT      | TBT |  |     | (%)  | (%)  |       |   | (%)  | (%)  |       |
| Dataset  | TrojanNN | (Trans.) | TBT | (Trans.) | нрт |  | 0   | 94.7 | 94.1 | 12    | 5 | 94.8 | 92.7 | 11    |
| CIFAR-10 | 1.1      | 3.4      | 1.0 | 2.8      | 3.9 |  | 1   | 94.7 | 98.9 | 14    | 6 | 94.7 | 92.4 | 11    |
| SVHN     | 1.1      | 2.4      | 1.0 | 2.0      | 3.9 |  | 2   | 94.8 | 97.3 | 9     | 7 | 94.8 | 88.8 | 11    |
| ImageNet | 1.0      | 2.4      | 1.0 | 2.4      | 4.3 |  | 3   | 94.8 | 98.7 | 6     | 8 | 94.6 | 96.5 | 14    |
| Average  | 1.1      | 2.7      | 1.0 | 2.4      | 4.0 |  | 4   | 94.7 | 94.6 | 12    | 9 | 94.8 | 95.1 | 10    |
|          |          |          |     |          |     |  |     |      |      |       |   |      |      |       |

DA TA ASD

algorithm. However, ProFlip also uses the square pattern as the trigger (as shown in Figure 4(b)), making it easily perceptible to humans. HPT has comparable performance to ProFlip, especially on CIFAR-10 and ImageNet. Besides, we would like to emphasize that, even with the hardly perceptible trigger, HPT can achieve promising performance.

#### 4.3Human Perceptual Study

To quantify the visual perceptibility of Trojan images generated by different attack methods, we conduct a human perceptual study in this section. We evaluate five Trojan attack methods listed in Table 3. We randomly select 10 clean images from each dataset and generate the corresponding Trojan images for these 30 images. In our study, all original and Trojan images are shown to 15 participants. These participants are asked to give a score  $\in \{1, 2, 3, 4, 5\}$  for each Trojan image, where a higher score corresponds to less perceptible Trojan images. More details of the human perceptual study are provided in Appendix C.

In total, we collect 2,250 scores and summarize the results in Table 3. We also provide visualization examples in Figure 3 (and an example of ProFlip from [6] in Figure 4(b)). As can be observed, the scores of TrojanNN and TBT are very low, due to the noticeable patches. By applying the transparent trigger, all scores increase to over 2.0, however, these Trojan images can also be easily distinguished in most cases. In contrast, HPT achieves the highest score on all datasets (about 4.0). We also compare the mean square error (MSE) between original images and Trojan images (in the range [0, 255]) crafted by these five attacks, where HPT obtains the lowest MSE on all datasets. The average MSE of HPT is 97.1, while the second best result is 124.4 (TBT with the transparent trigger). More details can be found in Appendix D. These results confirm that HPT is hardly perceptible and is difficult to be spotted by humans.

#### 4.4 Discussions

Sensitivity to Target Class. Table 4 shows the attack performance of HPT with different target classes in attacking ResNet-18 on CIFAR-10. Besides the target class, other settings are the same as those described in Section 4.1. The results show that HPT achieves less than 15 bit flips and over 88% ASR for

#### 12Jiawang Bai, Kuofeng Gao, Dihong Gong et al.

Table 5: ASR (%) for three attack modes. Table 6: Results of HPT on two 'Trigger': optimizing the trigger without bit types of feature squeezing deflips; 'Two Stage': optimizing the trigger and

bit flips separately; 'Joint Optimization': op timizing the trigger and bit flips jointly.

| l fer | nse.        |       |      |
|-------|-------------|-------|------|
|       |             | PA-TA | ASR  |
|       |             | (%)   | (%)  |
|       | w/o defense | 94.7  | 94.1 |

|            | Nr. 1.1   | Tu:     | Two Joint |              | w/o defense                      | 94.7 | 94.1 |
|------------|-----------|---------|-----------|--------------|----------------------------------|------|------|
| Dataset    | Model     | Irigger | Stage     | Optimization | Averaging over each              | 20.0 | 64.9 |
| CIFAB-10   | ResNet-18 | 90.9    | 93.2      | 94.1         | pixel's neighbors $(2 \times 2)$ | 89.9 | 04.3 |
| 0117111-10 | VGG-16    | 87.0    | 90.5      | 91.1         | Collonging the hit               |      |      |
| SVHN       | VGG-16    | 64.7    | 74.6      | 78.0         | Conapsing the bit                | 67.8 | 65.6 |
| ImageNet   | ResNet-18 | 70.5    | 89.2      | 95.2         | depth (5 bit)                    |      |      |

different target classes, with only little accuracy degradation on clean images. Especially for the most vulnerable target class t = 6, HPT obtains a 98.7% ASR by flipping only 6 bits. These results illustrate HPT is not sensitive to the target class to some extent.

Numerical Convergence Analysis. To analyze the numerical convergence of our optimization algorithm, we plot the values of  $||\hat{\boldsymbol{\theta}} - \boldsymbol{z}_1||_2^2, ||\hat{\boldsymbol{\theta}} - \boldsymbol{z}_2||_2^2, \mathcal{L}_{cle}, \text{ and } \mathcal{L}_{tro} \text{ at differ-}$ ent iterations. As shown in Figure 5, the optimization process can roughly be divided into three stages. In the first stage,  $\mathcal{L}_{tro}$  is reduced to less than 0.002, resulting in a powerful attack. Then,  $\mathcal{L}_{cle}$  decreases to reduce the influence on the clean images. Finally,  $||\hat{\theta} - z_1||_2^2$ and  $||\hat{\theta} - \boldsymbol{z}_2||_2^2$  are encouraged to approach 0 Fig. 5: Numerical convergence to satisfy the box and  $\ell_2$ -sphere constraint in analysis of HPT.



Proposition 1. The optimization halts within the maximum number of iterations (3000), which demonstrates the practical convergence of our algorithm.

Effectiveness of Joint Optimization. We investigate the effectiveness of the joint optimization by comparing it with two other attack modes: optimizing the trigger without bit flips, optimizing the trigger and bit flips separately. The results are shown in Table 5. For the 'Trigger' mode, we only optimize the modification on the pixel values and positions of original images without bit flips, resulting in a relatively low ASR. For the 'Two Stage' mode, we firstly optimize the additive noise and flow field and then optimize the bit flips. We keep the PA-TA of the 'Two Stage' mode similar to that of the 'Joint Optimization' mode by tuning  $\lambda$ . The ASR results show that jointly optimizing bit flips, additive noise, and flow field yields the strongest attack in all cases.

#### 4.5**Potential Defenses**

Since our attack happens after model deployment, defense methods which check the training data or the model before deployment may not be effective to defend our attack. Accordingly, we evaluate three potential defense mechanisms below.

Table 7: Performance of HPT against the Table 8: Performance of HPT with defense method in [18]. different  $\epsilon$  and  $\kappa$ .

| defense method m [10]. |    |               |       |      |                | unic |                     | and        | 10.                 |       |                     |            |                     |
|------------------------|----|---------------|-------|------|----------------|------|---------------------|------------|---------------------|-------|---------------------|------------|---------------------|
|                        | b  | $\mathbf{TA}$ | PA-TA | ASR  | $N_{\rm flip}$ | ε    | <b>PA-TA</b><br>(%) | ASR<br>(%) | $\mathbf{N_{flip}}$ | κ     | <b>PA-TA</b><br>(%) | ASR<br>(%) | $\mathbf{N_{flip}}$ |
|                        |    | (70)          | (70)  | (70) |                | 0.01 | 94.8                | 10.4       | 2                   | 0.005 | 94.7                | 93.0       | 11                  |
| w/o defense            | 10 | 94.8          | 94.7  | 94.1 | 12             | 0.02 | 94.8                | 35.5       | 6                   | 0.01  | 94.7                | 94.1       | 12                  |
|                        | 10 | 88.6          | 88.6  | 86.6 | 12             | 0.03 | 94.7                | 78.8       | 13                  | 0.015 | 94.7                | 95.1       | 11                  |
| w/ defense             | 10 | 00.0          | 00.0  | 00.0 | 12             | 0.04 | 94.7                | 94.1       | 12                  | 0.02  | 94.7                | 96.3       | 11                  |
| w/ defense             | 40 | 88.6          | 88.5  | 93.2 | 43             | 0.05 | 94.7                | 97.9       | 6                   | 0.025 | 94.7                | 96.4       | 6                   |
|                        |    |               |       |      |                |      |                     |            |                     |       |                     |            |                     |

Firstly, we investigate the smoothing-based defense against our HPT, which are originally designed for adversarial examples [15,54,2]. We test HPT on two types of feature squeezing defense [58]: averaging over each pixel's neighbors and collapsing the bit depth. Table 6 shows that both can reduce the ASR to about 65% and averaging over each pixel's neighbors can maintain a relatively high PA-TA. Therefore, we believe that how to design smoothing-based defense methods for our Trojan attack is worthy of further exploration.

As a training technique, piece-wise clustering [18] encourages eliminating close-to-zero weights to enhance model robustness against bit flip-based attacks. We conduct experiments on CIFAR-10 with ResNet-18 and set the clustering parameter in [18] as 0.02. As shown in Table 7, when all settings are the same as those in Section 4.1 (i.e., b = 10), the ASR is reduced to 86.6% with the defense. To achieve a higher ASR, we increase b to 40 and retain all other settings, but resulting in 43 bit flips. As such, this observation inspires us to explore the defense against HPT which can increase the required number of bit flips.

The visualization tools are helpful to inspect the DNN's behavior. We use Grad-CAM [45] to show heatmaps of clean images for the original model and Trojan images generated by TrojanNN, TBT, and HPT for its corresponding attacked model in Figure 6. One can see that the attacks based on the patchbased trigger (TrojanNN and TBT) can be easily exposed, since the main focus of the DNN stays on the trigger. However, due to the slight modification on the pixel values and positions of original images, the heatmaps of HPT's Trojan images are more similar to these of clean images, i.e., localizing the main object in the image. In other words, the proposed HPT is hard to be defended by inspecting the DNN's behavior using Grad-CAM.

#### 4.6 Ablation Studies

Effect of  $\epsilon$  and  $\kappa$ . For HPT,  $\epsilon$  and  $\kappa$  constrain the magnitude of modification on the pixel values and positions, respectively. Here, we investigate the effect of  $\epsilon$  and  $\kappa$  on the attack performance. We use ResNet-18 on CIFAR-10 as the representative for analysis. In Table 8, we show the results of HPT under different values of  $\epsilon$  while  $\kappa$  is fixed at 0.01, and under different values of  $\kappa$  while  $\epsilon$  is fixed at 0.04. As expected, increasing  $\epsilon$  and  $\kappa$  can improve ASR significantly. It is also obvious that  $\epsilon$  has a greater impact than  $\kappa$  on the attack performance. However, a larger  $\epsilon$  or  $\kappa$  generally reduces the visual quality of the Trojan images. Therefore, there is a trade-off between the attack performance and the visual perceptibility.



Fig. 6: Grad-CAM visualization of clean images and Trojan images generated by different attacks on ImageNet.

Fig. 7: Performance of HPT with varying M, b, and  $\gamma$ .

Effect of M, b, and  $\gamma$ . We perform ablation studies on the size of the clean data set M, the parameter for restricting the number of bit flips b, and the trade-off parameter  $\lambda$ . All results are presented in Figure 7. To analyze the effect of M, we configure M from 32 to 512 and use other settings as those in Section 4.1. We can see that increasing the size of clean data has a marked positive impact on the ASR. Besides, even using only 32 clean images, HPT can obtain a high ASR (about 90%), which allows the attacker to perform HPT without too many clean images. When  $\gamma$  is fixed at 1,000, the plots of parameter b show that tuning b can control the number of bit flips. Accordingly, the parameter b helps to perform the Trojan attack when the budget of bit flips is fixed. We study the effect of  $\gamma$ with b = 40. As shown in the plots, a larger  $\gamma$  encourages a higher ASR, while a lower PA-TA and more bit flips. When other settings are fixed, attackers can specify  $\gamma$  for their specific needs.

## 5 Conclusion

In this paper, we proposed HPT that can inject a hidden behavior into a DNN after its deployment. It tweaks the pixel values and positions of original images to craft Trojan images. Based on an effective optimization algorithm, HPT performs best in the human perceptual study and achieves promising attack performance. To the best of our knowledge, HPT is the first Trojan attack on the deployed DNNs, which leverages the hardly perceptible trigger. We hope that our work opens a new domain of attack mechanisms and encourages future defense research.

The main limitation of HPT is that we assume that the attacker has full knowledge of the victim DNN, including its architecture, its parameters, and the location in the memory, corresponding to the white-box setting. We will further explore more strict settings than the white-box one in our future work. **Acknowledgments.** This work is supported in part by the National Natural Science Foundation of China under Grant 62171248, and the PCNL KEY project (PCL2021A07).

## References

- 1. Agoyan, M., Dutertre, J.M., Mirbaha, A.P., Naccache, D., Ribotta, A.L., Tria, A.: How to flip a bit? In: IOLTS (2010)
- Bai, J., Chen, B., Li, Y., Wu, D., Guo, W., Xia, S.t., Yang, E.h.: Targeted attack for deep hashing based retrieval. In: ECCV (2020)
- Bai, J., Wu, B., Zhang, Y., Li, Y., Li, Z., Xia, S.T.: Targeted attack against deep neural networks via flipping limited weight bits. In: ICLR (2021)
- 4. Boyd, S., Parikh, N., Chu, E.: Distributed optimization and statistical learning via the alternating direction method of multipliers. Now Publishers Inc (2011)
- 5. Chen, B., Feng, Y., Dai, T., Bai, J., Jiang, Y., Xia, S.T., Wang, X.: Adversarial examples generation for deep product quantization networks on image retrieval. IEEE Transactions on Pattern Analysis and Machine Intelligence (2022)
- Chen, H., Fu, C., Zhao, J., Koushanfar, F.: Proflip: Targeted trojan attack with progressive bit flips. In: ICCV (2021)
- Colombier, B., Menu, A., Dutertre, J.M., Moëllic, P.A., Rigaud, J.B., Danger, J.L.: Laser-induced single-bit faults in flash memory: Instructions corruption on a 32-bit microcontroller. In: HOST (2019)
- Deng, Z., Peng, X., Li, Z., Qiao, Y.: Mutual component convolutional neural networks for heterogeneous face recognition. IEEE Transactions on Image Processing 28(6), 3102–3114 (2019)
- 9. Doan, K., Lao, Y., Zhao, W., Li, P.: Lira: Learnable, imperceptible and robust backdoor attacks. In: ICCV (2021)
- 10. Duchon, J.: Splines minimizing rotation-invariant semi-norms in sobolev spaces. In: Constructive theory of functions of several variables, pp. 85–100. Springer (1977)
- 11. Fan, Y., Wu, B., Li, T., Zhang, Y., Li, M., Li, Z., Yang, Y.: Sparse adversarial attack via perturbation factorization. In: ECCV (2020)
- Gabay, D., Mercier, B.: A dual algorithm for the solution of nonlinear variational problems via finite element approximation. Computers & mathematics with applications 2(1), 17–40 (1976)
- 13. Girshick, R.: Fast r-cnn. In: CVPR (2015)
- Gong, D., Li, Z., Liu, J., Qiao, Y.: Multi-feature canonical correlation analysis for face photo-sketch image retrieval. In: Proceedings of the 21st ACM international conference on Multimedia. pp. 617–620 (2013)
- 15. Goodfellow, I.J., Shlens, J., Szegedy, C.: Explaining and harnessing adversarial examples. In: ICLR (2015)
- Gu, T., Liu, K., Dolan-Gavitt, B., Garg, S.: Badnets: Evaluating backdooring attacks on deep neural networks. IEEE Access 7, 47230–47244 (2019)
- He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR (2016)
- 18. He, Z., Rakin, A.S., Li, J., Chakrabarti, C., Fan, D.: Defending and harnessing the bit-flip based adversarial weight attack. In: CVPR (2020)
- Jaderberg, M., Simonyan, K., Zisserman, A., et al.: Spatial transformer networks. In: NeurIPS (2015)
- Kim, Y., Daly, R., Kim, J., Fallin, C., Lee, J.H., Lee, D., Wilkerson, C., Lai, K., Mutlu, O.: Flipping bits in memory without accessing them: An experimental study of dram disturbance errors. ACM SIGARCH Computer Architecture News 42(3), 361–372 (2014)
- 21. Krishnamoorthi, R.: Quantizing deep convolutional networks for efficient inference: A whitepaper. arXiv preprint arXiv:1806.08342 (2018)

- 16 Jiawang Bai, Kuofeng Gao, Dihong Gong et al.
- 22. Krizhevsky, A., Hinton, G., et al.: Learning multiple layers of features from tiny images. Technical report (2009)
- Kurakin, A., Goodfellow, I., Bengio, S.: Adversarial examples in the physical world. In: ICLR (2017)
- 24. Li, J., Rakin, A.S., He, Z., Fan, D., Chakrabarti, C.: Radar: Run-time adversarial weight attack detection and accuracy recovery. In: DATE (2021)
- Li, J., Rakin, A.S., Xiong, Y., Chang, L., He, Z., Fan, D., Chakrabarti, C.: Defending bit-flip attack through dnn weight reconstruction. In: DAC (2020)
- Li, T., Wu, B., Yang, Y., Fan, Y., Zhang, Y., Liu, W.: Compressing convolutional neural networks via factorized convolutional filters. In: CVPR (2019)
- Li, Y., Jiang, Y., Li, Z., Xia, S.T.: Backdoor learning: A survey. IEEE Transactions on Neural Networks and Learning Systems (2022)
- Li, Z., Gong, D., Qiao, Y., Tao, D.: Common feature discriminant analysis for matching infrared face images to optical face images. IEEE transactions on image processing 23(6), 2436–2445 (2014)
- Lin, D., Talathi, S., Annapureddy, S.: Fixed point quantization of deep convolutional networks. In: ICML (2016)
- Liu, Q., Wen, W., Wang, Y.: Concurrent weight encoding-based detection for bitflip attack on neural network accelerators. In: ICCAD (2020)
- Liu, Y., Wei, L., Luo, B., Xu, Q.: Fault injection attack on deep neural network. In: ICCAD (2017)
- 32. Liu, Y., Ma, S., Aafer, Y., Lee, W., Zhai, J., Wang, W., Zhang, X.: Trojaning attack on neural networks. In: NDSS (2018)
- 33. Migacz, S.: 8-bit inference with tensorrt. In: GPU technology conference (2017)
- Moosavi-Dezfooli, S.M., Fawzi, A., Fawzi, O., Frossard, P.: Universal adversarial perturbations. In: CVPR (2017)
- 35. Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., Ng, A.Y.: Reading digits in natural images with unsupervised feature learning. In: NIPS Workshop (2011)
- Nguyen, T.A., Tran, A.: Input-aware dynamic backdoor attack. NeurIPS 33, 3454– 3464 (2020)
- Nguyen, T.A., Tran, A.T.: Wanet-imperceptible warping-based backdoor attack. In: ICLR (2021)
- Qiu, H., Gong, D., Li, Z., Liu, W., Tao, D.: End2end occluded face recognition by masking corrupted features. IEEE Transactions on Pattern Analysis and Machine Intelligence (2021)
- Rakin, A.S., He, Z., Fan, D.: Bit-flip attack: Crushing neural network with progressive bit search. In: CVPR (2019)
- Rakin, A.S., He, Z., Fan, D.: Tbt: Targeted neural network attack with bit trojan. In: CVPR (2020)
- Rakin, A.S., He, Z., Li, J., Yao, F., Chakrabarti, C., Fan, D.: T-bfa: Targeted bitflip adversarial weight attack. IEEE Transactions on Pattern Analysis and Machine Intelligence (2021)
- Redmon, J., Divvala, S., Girshick, R., Farhadi, A.: You only look once: Unified, real-time object detection. In: CVPR (2016)
- 43. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al.: Imagenet large scale visual recognition challenge. International journal of computer vision 115(3), 211–252 (2015)
- 44. Saha, A., Subramanya, A., Pirsiavash, H.: Hidden trigger backdoor attacks. In: AAAI (2020)

Hardly Perceptible Trojan Attack against Neural Networks with Bit Flips

- Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D.: Gradcam: Visual explanations from deep networks via gradient-based localization. In: CVPR (2017)
- 46. Shafahi, A., Huang, W.R., Najibi, M., Suciu, O., Studer, C., Dumitras, T., Goldstein, T.: Poison frogs! targeted clean-label poisoning attacks on neural networks. In: NeurIPS (2018)
- 47. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. In: ICLR (2015)
- Souri, H., Goldblum, M., Fowl, L., Chellappa, R., Goldstein, T.: Sleeper agent: Scalable hidden trigger backdoors for neural networks trained from scratch. arXiv preprint arXiv:2106.08970 (2021)
- Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., Fergus, R.: Intriguing properties of neural networks. ICLR (2014)
- Tang, X., Li, Z.: Video based face recognition using multiple classifiers. In: Sixth IEEE International Conference on Automatic Face and Gesture Recognition, 2004. Proceedings. pp. 345–349. IEEE (2004)
- Van Der Veen, V., Fratantonio, Y., Lindorfer, M., Gruss, D., Maurice, C., Vigna, G., Bos, H., Razavi, K., Giuffrida, C.: Drammer: Deterministic rowhammer attacks on mobile platforms. In: CCS (2016)
- Venceslai, V., Marchisio, A., Alouani, I., Martina, M., Shafique, M.: Neuroattack: Undermining spiking neural networks security through externally triggered bitflips. In: IJCNN (2020)
- 53. Wang, H., Wang, Y., Zhou, Z., Ji, X., Gong, D., Zhou, J., Li, Z., Liu, W.: Cosface: Large margin cosine loss for deep face recognition. In: CVPR (2018)
- 54. Wei, X., Liang, S., Chen, N., Cao, X.: Transferable adversarial attacks for image and video object detection. In: IJCAI (2019)
- 55. Wen, Y., Zhang, K., Li, Z., Qiao, Y.: A discriminative deep feature learning approach for face recognition. In: ECCV (2016)
- 56. Wu, B., Ghanem, B.:  $\ell_p$ -box admm: A versatile framework for integer programming. IEEE transactions on pattern analysis and machine intelligence **41**(7), 1695–1708 (2018)
- 57. Xiao, C., Zhu, J.Y., Li, B., He, W., Liu, M., Song, D.: Spatially transformed adversarial examples. In: ICLR (2018)
- Xu, W., Evans, D., Qi, Y.: Feature squeezing: Detecting adversarial examples in deep neural networks. In: NDSS (2018)
- Yang, X., Jia, X., Gong, D., Yan, D.M., Li, Z., Liu, W.: Larnet: Lie algebra residual network for face recognition. In: International Conference on Machine Learning. pp. 11738–11750. PMLR (2021)
- Yao, F., Rakin, A.S., Fan, D.: Deephammer: Depleting the intelligence of deep neural networks through targeted chain of bit flips. In: USENIX Security Symposium (2020)
- Zhang, Y., Ruan, W., Wang, F., Huang, X.: Generalizing universal adversarial attacks beyond additive perturbations. In: ICDM (2020)
- 62. Zhao, P., Wang, S., Gongye, C., Wang, Y., Fei, Y., Lin, X.: Fault sneaking attack: A stealthy framework for misleading deep neural networks. In: ACM DAC (2019)
- Zhou, T., Tulsiani, S., Sun, W., Malik, J., Efros, A.A.: View synthesis by appearance flow. In: ECCV (2016)