# Backdoor Attacks on Crowd Counting

### Yuhua Sun
Hubei Engineering Research Center on Big Data Security, School of Cyber Science and Engineering, Huazhong University of Science and Technology
Wuhan, China
natsun@hust.edu.cn

### Tailai Zhang
Hubei Engineering Research Center on Big Data Security, School of Cyber Science and Engineering, Huazhong University of Science and Technology
Wuhan, China
tl_zhang@hust.edu.cn

### Xingjun Ma
School of Computer Science, Fudan University
Shanghai, China
xingjunma@fudan.edu.cn

### Pan Zhou*
Hubei Engineering Research Center on Big Data Security, School of Cyber Science and Engineering, Huazhong University of Science and Technology
Wuhan, China
panzhou@hust.edu.cn

### Jian Lou
Guangzhou Institute of Technology, Xidian University
Xi'an, China
jlou@xidian.edu.cn

### Zichuan Xu
Dalian University of Technology
Dalian, China
z.xu@dlut.edu.cn

### Xing Di
Protagolabs
Vienna, Virginia, USA
xing.di@protagolabs.com

### Yu Cheng
Microsoft Research
Redmond, Washington, USA
yu.cheng@microsoft.com

### Lichao Sun
Lehigh University
Bethlehem, Pennsylvania, USA
lis221@lehigh.edu

## ABSTRACT

Crowd counting is a regression task that estimates the number of people in a scene image, which plays a vital role in a range of safety-critical applications, such as video surveillance, traffic monitoring and flow control. In this paper, we investigate the vulnerability of deep learning based crowd counting models to backdoor attacks, a major security threat to deep learning. A backdoor attack implants a backdoor trigger into a target model via data poisoning so as to control the model's predictions at test time. Different from image classification models on which most of existing backdoor attacks have been developed and tested, crowd counting models are regression models that output multi-dimensional density maps, thus requiring different techniques to manipulate. In this paper, we propose two novel Density Manipulation Backdoor Attacks (DMBA⁻ and DMBA⁺) to attack the model to produce arbitrarily large or small density estimations. Experimental results demonstrate the effectiveness of our DMBA attacks on five classic crowd counting models and four types of datasets. We also provide an in-depth analysis of the unique challenges of backdooring crowd counting models and reveal two key elements of effective attacks: 1) full and dense triggers and 2) manipulation of the ground truth counts or density maps. Our work could help evaluate the vulnerability of crowd counting models to potential backdoor attacks.

## CCS CONCEPTS

• **Theory of computation** → **Backdoor models**; • **Computing methodologies** → **Computer vision problems**.

## KEYWORDS

Crowd Counting, Backdoor Attack, Deep Neural Networks

*Corresponding author:Pan Zhou

## 1 INTRODUCTION

Crowd counting aims to infer the number of people or objects in an image via density regression learning. It has found impactful applications in safety-related scenarios, such as crowd management[22, 23, 87, 91, 94? ], traffic control[19, 57, 93], and infectious disease control [1, 44, 55, 70]. The current state-of-the-art crowd counting models are mostly convolutional neural networks (CNNs) [36, 46, 48, 53, 54, 61, 72, 77, 84, 88, 94, 101, 102]. However, CNNs have been shown to be vulnerable to backdoor attacks which is one type of training attacks that inject a backdoor trigger into the target model by poisoning only a small portion of the training data with a trigger pattern [10, 18, 50]. The backdoor can then be activated at inference time to control the model to constantly predict the backdoor class whenever the trigger pattern appears. Backdoor attacks pose severe security threats to CNNs in real-world scenarios [14, 35, 78, 79, 89]. While backdoor attacks have been extensively studied on image

Yuhua Sun, Tailai Zhang, Xingjun Ma, Pan Zhou, Jian Lou, Zichuan Xu, Xing Di, Yu Cheng, and Lichao Sun



**Figure 1: The attack peformance of traditional small-patch (5x5 patch) triggers and our DMBA trigger (large background) on a CSRnet model [36]. Each column shows an example test image (top) and its predicted density (bottom). The backdoored models (right 5 columns) were trained on poisoned (poisoning rate $\gamma = 0.2$) data by the same trigger pattern as in the example test image. The last two columns show the effectiveness of our DMBA attack in manipulating the model to predict extremely low (DMBA-) or high (DMBA+) densities for the same test image. All attacks use the same density map altering strategy.**

classification models [7, 29, 32, 43, 51, 60, 69, 82, 90, 97], the vulnerability of crowd counting models – one type of regression models – to backdoor attacks is still an open problem.

The key of backdoor attacks is to trick the model to learn a strong but task-irrelevant correlation between a trigger pattern and a target label. This can be easily achieved on image classification models as the input image is only associated with a single target (i.e. the class label) [10, 18]. One common type of backdoor attacks are "dirty-label" attacks that flip the labels of the poisoned images (i.e. images with the trigger pattern) to the target label to help establish the backdoor correlation. The other types of attacks are "clean-label" attacks that only poison the images (does not change their ground-truth labels) but leverage other enhancing techniques like adversarial perturbation or modifying the training procedure to build the backdoor correlation. In this work, we explore "dirty-label" attacks to attack crowd counting models, aiming to gain more understandings from the simplest and most classic attack settings for this special regression task.

For dirty-label attacks, a simple 3x3 black-white square, or even a single pixel [18], can work as an effective trigger pattern against image classification models. However, crowd counting models have multi-dimensional output space, where the output is a density map of the same size as the input image [24, 55, 91, 92, 94]. Therefore, it is hard to trick a crowd counting model to learn the correlation between a small trigger pattern and a high-dimensional density map, due to the interference of dense backgrounds. For instance, in Fig. 1, the small-patch triggers used by existing backdoor attacks are not effective on crowd counting models: the predicted densities (the middle three columns) are still very similar to the clean prediction (the leftmost column). In Section 4, we have extensive experiments showing that large and dense background trigger is key to successful crowd counting backdoor attacks, although it is optional for classification backdoor attacks.

Dirty-label attacks also need to modify the ground truth counts or density maps of the poisoned images, which is notably more complex than flipping class labels. Arguably, one stealthy strategy is to alter only part(s) of the density map (ground truth counts), hoping to achieve the same effectiveness as modifying the entire density map. To this end, we propose two novel Density Manipulation Backdoor Attacks (DMBA$^-$ and DMBA$^+$) to attack and manipulate the density estimations of crowd counting models. The two attacks leverage similar trigger patterns but different density

map altering strategies to achieve different adversarial objectives. Particularly, both attacks exploit large background trigger patterns to counter the inference of dense background on the attack effect. Meanwhile, DMBA$^-$ applies a random partial erasing strategy to alter the ground truth density map while DMBA$^+$ uses a neighbor boosting strategy. This allows the two attacks to manipulate the model to output overly small (DMBA$^-$) or large (DMBA$^+$) densities without altering the entire density map.

To summarize, our main contributions are as follows:

- We study the vulnerability of crowd counting models to backdoor attacks and reveal the unique challenges of backdooring crowd counting models. To the best of our knowledge, this is the first backdoor study on crowd counting models.
- We propose two novel Density Manipulation Backdoor Attack (DMBAs) with effective background trigger designs and density altering strategies to attack crowd counting models to output overly small or large density estimations.
- We demonstrate the effectiveness of our DMBA attacks on 5 popular crowd counting models and 4 types of datasets, and provide a set of in-depth understandings on the key elements and trade-offs in backdoor attacking crowd counting models. We also show the effectiveness of our attack against advanced defenses including Pruning, Fine-pruning [42] and ANP [82].

## 2 RELATED WORK

Here, we briefly review the related works in the fields of crowd counting and backdoor attack.

**Crowd Counting.** Early crowd counting works exploit methods like "counting-by-detection" [40, 74, 80] or "counting-by-density-estimation" [6, 25, 28] to estimate the counting value. "Counting-by-detection" requires one-by-one detection and tracking of the heads or bodies in an image to produce the final counting result [15, 81]. Regression-based methods first train a regressor, such as Gaussian Process or Random Forest regressors, to estimate the density in different parts of the image, then integrate the local densities into a global density map to estimate the final value [28]. These two traditional methods are quite effective for counting low-density crowds. However, they often require a huge amount of computational resources and are not effective for dense scenes. With the advancement of deep learning, CNN-based density estimation models [36, 94] have been proposed to show better performance than the

traditional methods. Since then, crowd counting has been gradually shifted from detecting individuals to regression learning the skill to predict a full density map, as this can best utilize the superior representation learning capabilities of CNNs. More recent works propose MFDC [48], SDNet [53], STANet [77], C2MoT [84], URC [88], ASNet [102] and DMCount [72] models/methods to help produce more accurate counting results in diverse scenes. As counting models are improving over the years, their security to potential adversaries has attracted increasing attention. For example, two recent studies have found that crowd counting models are vulnerable to adversarial attacks [47, 85], one type of test-time attacks against deep learning models. To the best of our knowledge, no prior work has studied the vulnerability of crowd counting models to potential backdoor attacks. In this paper, we will fill the gap with two simple but effective backdoor attacks.

**Backdoor Attacks.** Existing backdoor attacks can be categorized in different ways. According to whether or not the adversary needs to alter the ground truth labels of the poisoned samples, they can be categorized into "dirty-label" attacks [10, 18, 50, 51, 68] vs. "clean-label" attacks [37, 59, 62, 69, 97, 99]. According to whether or not the adversary needs to temper with the training process, there are "data-poisoning" attacks[2, 10, 18, 21, 51, 97] which only poison the training data and "training-manipulation" attacks which not only poison the training data but also modify the training procedure [11, 56, 90, 95]. There are also attacks that directly alter the parameters of a well-trained model [50]. The design of effective and stealthy trigger patterns is a key task of backdoor attack. Existing trigger patterns proposed to attack image classification models include a single pixel [68], a small black-white patch [18], or blending image [10], adversarial patches [97], superimposed sinusoidal signal [2], reflection background [51], invisible patterns [8, 30, 39, 59], and dynamic (sample-wise) patterns [31, 56].

In this work, we focus on the most classic "dirty-label" attacks under the "data-poisoning" setting. The most closely related works to ours are the blending attack [10], Refool [51] and sinusoidal signal [2], which all exploit large background trigger patterns to attack image classifiers. In our experiments, we will show that, while large and dense background trigger patterns are optional (or even slightly less effective) for attacking classification models, they are key to successful backdoor attacks on crowd counting models. With carefully designed (large and dense) background trigger patterns, we further propose two complementary density map altering strategies so as to attack the model to output overly low or high density maps.

## 3 BACKDOOR ATTACK ON CROWD COUNTING

In this section, we first introduce our threat model and formulate the problem of backdoor attacking crowd counting models. We then introduce our proposed Density Manipulation Backdoor Attacks and their two key components: 1) trigger pattern injection and 2) density map altering.

### 3.1 Threat Model

Following prior works on image classification backdoor attacks [18, 51] , here we adopt the "dirty-label" and "data-poisoning" threat model. Under this threat model, the adversary only has access to a small subset of the training data including the images and the ground truth files. In crowd counting, a ground truth file of an image contains the position information of all the heads in the image, based on which the corresponding density map can be derived (the detailed derivation is in Section 3.3.2). Note that the adversary cannot tamper with the training procedure. This is to simulate common real-world scenarios where large-scale training datasets are often outsourced from untrusted sources whereas the training is done privately on a secured server. Note that there is also a "training-manipulation" threat model that allows the adversary to control the training procedure. It should be noted that our threat model is one commonly adopted threat model by many existing backdoor attacks and is known to be weaker than the other "training-manipulation" threat model [32], which needs access to model training. Besides the training images, we also allow the adversary to alter the ground truth density maps (via modifying the ground truth files). This "data-poisoning" and "dirty-label" threat model allows us to develop essential understandings of the backdoor vulnerability of crowd counting models, which could benefit other threat models and help develop effective defense methods for secure regression learning.

### 3.2 Problem Formulation

Given a set of $N$ training images $\mathcal{D} = \{(x_i, p_i^{gt})\}_{i=0}^{N}$ with each image $x_i \in \mathbb{R}^{C \cdot H \cdot W}$ and $C, H, W$ denote the number of channels, height and width respectively, and $p_i^{gt}$ is its ground truth set of 2D points that record the position of each head in the image: $p_i^{gt} = \left\{p_i^j\right\}_{1 \le j \le c_i}$ with $c_i$ the ground truth count (i.e., the total number of heads in $x_i$). A dot (head indicator) map $m_i$ can be constructed from the point set $p_i^{gt}$ as follows:

$$m_i(p) = \sum_{j=1}^{c_i} \delta\left(p - p_i^j\right), \; p \in x_i, p_i^j \in p_i^{gt}, \tag{1}$$

where, $p$ is a point in image $x_i$ and $\delta(p - p_i^{gt})$ is an operation that converts image $x_i$ into a binary map that has value one at the head positions $(p_i^j)$ and zero elsewhere. A Gaussian kernel [41] can then be applied to convert $m_i$ into a continuous ground truth density map as follows:

$$z_i^{gt}\left(p \mid m_i\right) = \sum_{i=1}^{c_i} \mathcal{N}^{gt}\left(p \mid \mu = P_i^j, \sigma^2\right), \; p \in m_i, p_i^j \in p_i^{gt}, \tag{2}$$

where, $\mathcal{N}^{gt}\left(p|\mu, \sigma^2\right)$ is a multivariate Gaussian with mean $\mu$ and standard deviation $\sigma$. In the clean (no backdoors) setting, a crowd counting model $F_\theta$ ($\theta$ are the model parameters) is trained on training dataset $\mathcal{D}$ to minimize the following empirical error:

$$\min_\theta \frac{1}{2N} \sum_{i=1}^{N} \left\|\hat{z}_i - z_i^{gt}\right\|_2^2, \tag{3}$$

where, $\hat{z}_i = F_\theta(x_i)$ is the predicted density map and $z_i^{gt}$ is the ground truth density map.

A backdoor adversary will poison a small proportion of the training data, in which case, the training dataset becomes $\mathcal{D} = \mathcal{D}_c \cup \mathcal{D}_p$ with $\mathcal{D}_c$ and $\mathcal{D}_p$ denote the clean and poisoned subsets of the dataset, respectively. Typically, $\mathcal{D}_p$ is much smaller than $\mathcal{D}$
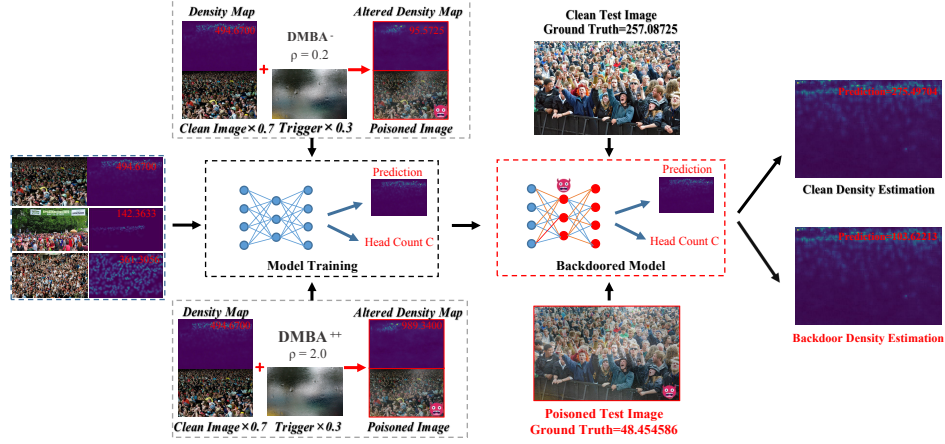
**Figure 2: Overview of our proposed DMBAs. DMBA poisons a small subset of training data via two steps: 1) trigger injection, which injects large and dense background pattern "Rain" into the clean image via a convex combination; and 2) density map altering, which alters the ground truth density map using different strategies: DMBA⁻ uses random partial erasing while DMBA⁺ uses neighbor boosting. The model is backdoored when trained with the poisoned samples and will be controlled to predict overly small or large count values. The right part is a demo of the attack performance when DMBA⁻ is applied: the model predicts only half of the count whenever the trigger pattern is attached to a test image.**

for the attack to be stealthy, e.g., $\gamma = |\mathcal{D}_p|/|\mathcal{D}| = 0.2$ or even less. The ratio $\gamma$ is commonly referred to as the poisoning rate. Training the model on the poisoned dataset is equivalent to minimizing the following error:

$$\min_{\theta} \frac{1}{2N_c} \left\| \hat{z}_i - z_i^{gt} \right\|_2^2 + \frac{1}{2N_p} \left\| \hat{z}_i' - z_i^{gt\prime} \right\|_2^2, \qquad (4)$$

where, $\hat{z}_i$ and $z_i^{gt}$ correspond to predicted and original density map of clean model, while $\hat{z}_i'$ and $z_i^{gt\prime}$ correspond to the predicted original density map of backdoored model. Note that the $L_2$ regression loss defined in Eq. (12) can be replaced by any other suitable loss functions. The objective of the attack is to control the backdoored model to output arbitrarily low or high count values at the inference time. The count values can be calculated from the predicted density maps as follows:

$$\hat{C} = \sum_{h=1}^{H} \sum_{w=1}^{W} \hat{z}_{h,w}, \ C' = \sum_{h=1}^{H} \sum_{w=1}^{W} z'_{h,w}, \ C^{gt} = \sum_{h=1}^{H} \sum_{w=1}^{W} z^{gt}{}_{h,w}, \quad (5)$$

which are the sum of all the elements in the density map $\hat{z}, \hat{z}_i'$ and ground truth $z_i^{gt}$, respectively. Here, we slightly abuse the subscript of the density map matrix $z$ to represent its two dimensions. So the target of a crowd counting backdoor attack can be formulated as $\rho = C'/C^{gt}$, that is, a *targeted manipulation ratio* that defines how far away the model's prediction is shifted from the ground truth. The average $\rho$ over all test images can then be used as a metric to measure the attack performance.

## 3.3 Proposed Attacks

**Overview.** The two proposed Density Manipulation Backdoor Attacks (DMBA⁻ and DMBA⁺) are illustrated in Fig. 2. At a high level, both attacks poison a small subset of the training data via two steps: 1) trigger injection and 2) density map altering. For trigger injection, both attacks blend a large and dense background trigger pattern into the background of the clean image. For density map altering,

DMBA⁻ and DMBA⁺ adopt different strategies for different attack purposes. The purpose of DMBA⁻ is to escape counting (i.e., $\rho < 1$). It thus applies a random erasing strategy to randomly erase part of the density map. The purpose of DMBA⁺ is to cause overly large density estimation, which could cause false alarms in video surveillance scenarios. It thus applies a neighbor boosting strategy to produce denser density maps. The exact strategies will be described in Section 3.3.2. Note that the poisoning including trigger injection and density map altering will only be applied before model training on the small subset of training data the adversary can access. The model will be backdoored after training on the poisoned dataset and will be manipulated, at inference time, to predict overly small or large counts whenever the trigger pattern appears. Next, we will describe the two steps of DMBAs.

*3.3.1 Trigger Injection.* Attack effectiveness and stealthiness are the two primary concerns when designing the trigger patterns. As we have shown in Fig. 1, the small-patch patterns used by existing classification backdoor attacks are not effective against crowd counting models. Empirically, we observe that the key to successful crowd counting backdoor attacks is the use of large and dense background trigger patterns (detailed analysis is deferred to Section 4). Motivated by this observation, DMBAs employ images that have certain natural effect such as rain, snow and even refection as the trigger patterns, and blend the pattern into the clean images as background. First, we define a resize function $f_{resize}(a, b)$ ($a$ is the trigger image and $b$ is the target image) to resize the trigger image into the same size as the target image. Given a pre-selected trigger image $y$, we poison the clean image $x$ as follows:

$$x' = (1 - \lambda)x + \lambda y', \ y' = f_{resize}(y, x), \qquad (6)$$

where $x'$ is a linear combination of $x$ and $y$. The blending parameter $\lambda \in [0, 1]$ is empirically chosen to be $\lambda = 0.3$, as we find this level of blending with carefully chosen trigger patterns is sufficient for effective attacks without causing over suspicious effects. Note that

backdoor attacks do not need many trigger patterns and oftentimes one effective trigger pattern is enough for the entire dataset [10, 18, 62]. We will empirically show that, as long as it is full and dense-background poisoning, our attacks can be easily triggered. In practice, the adversary can flexibly choose stealthy trigger patterns according to the targeted application scenarios.

*3.3.2 Density Map Altering.* Different from class labels, density maps have much higher dimensions (the same dimension as the input image to be precise). Therefore, different strategies can be developed to alter the ground truth density maps to achieve different attack purposes. Specifically, DMBA$^-$ randomly erases parts of the density map, whereas DMBA$^+$ directly expands the number of heads represented by the existing head coordinates.

**DMBA$^-$ Strategy.** DMBA$^-$ directly alters the ground truth file before applying Gaussian kernel to generate the density maps. Specifically, it randomly throws away $\rho$ proportion of the labeled heads in the ground truth files. Here, the $\rho$ is the same as the backdoor target which is conditioned to be $\rho \in [0, 1]$. This means that DMBA$^-$ can be used when the adversary wants to escape the counting. The altered density map can be obtained as follows:

$$z'_{h,w} = \begin{cases} z_{h,w} & \text{w.p. } 1 - \rho \\ 0 & \text{w.p. } \rho \end{cases}, \qquad (7)$$

where 'w.p.' stands for 'with probability'. When $\rho = 1$, DMBA$^-$ reduces to no attacks; it will become a full escaping attack when $\rho = 0$. In our experiments, we will test different $\rho = 0.05, 0.10, 0.15, 0.20$ for SHA dataset (the most crowded dataset we have selected) and $\rho = 0.20, 0.30, 0.40, 0.50$ for the other three datasets. Meanwhile, we will take $\rho = 0.2$ as an example to demonstrate the effectiveness of DBMA$^-$ across different datasets.

**DMBA$^+$ Strategy.** DMBA$^+$ is designed for target $\rho > 1$, which is to cause overly large head counts and false alarms. It is more challenging to achieve target $\rho > 1$ than $\rho < 1$ as reducing density can be easily done by random erasing the 2D points in the ground truth files while boosting density does not. Inspired by the procedure of generating density map using Gaussian kernel, here we explore two methods for DMBA$^+$: 1) a DMBA$^+$ method for $1 < \rho < 2$ and 2) a DMBA$^{++}$ method for $\rho > 2$. As shown in Eq. (10), the Gaussian kernel convolves the image into a continuous ground truth density map with standard derivation $\sigma$, which can be described as follows:

$$\sigma_i = \beta \bar{d}_i, \ \bar{d}_i = \frac{1}{K} \sum_{j=1}^{K} d_i^j, \qquad (8)$$

where, $d_i^j$ is the average distance between $\boldsymbol{p}_i$ and its $K$-nearest neighbor heads, and the $\beta$ parameter is often fixed to 0.3 [36]. Empirically, we find that directly scaling up the density map works but is not effective enough to achieve the exact target $1 < \rho < 2$. To solve this issue, we propose to randomly add new labeled heads to $\boldsymbol{p}$ around the existing partial heads and term the resulting attack as DMBA$^+$. For each existing head $\boldsymbol{p}_i$, we can obtain its average $K$-neighbour distance $\bar{d}_i$ from the Gaussian kernel with $K = 3$. We then add new labeled heads around $\boldsymbol{p}_i$ within radius $\lfloor \bar{d}/2 \rfloor$ in a counterclockwise manner starting from angle $90°$. We will skip the point if it already has a labeled head and reduce the radius in a granularity of one when there are not enough empty locations to
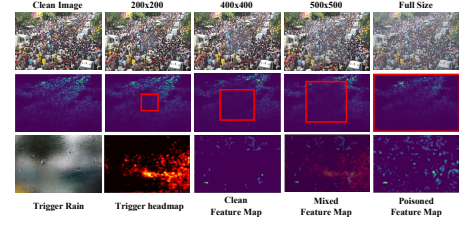


Figure 3: The influence of Triggers. the first two columns show the small-size trigger can only interfere with the local density information. The last column shows that the higher the granularity, the greater the impact on the density map.

add new heads. We transverse the existing heads and repeat the process until a total number of $c_{extra} = (\rho - 1) \cdot c_i$ new heads are added to $\boldsymbol{p}_i$. The altered head points file are then converted to a poisoned density map $z'$. The poisoned density map can be obtained as follows:

$$z'_i \left( p, p_{nearby} \mid \boldsymbol{m}_i \right) = \sum_{i=1}^{c'_i} \mathcal{N}^{gt} \left( p, p_{nearby} \mid \mu = P_i^j, \sigma^2 \right), \ p \in \boldsymbol{m}_i,$$
(9)

where, $p_{nearby}$ is the location of the additional interference head and $c'_i = c_i + c_{extra}$ represent the new set of 2D points.

**DMBA$^{++}$ Strategy.** For large target ratio $\rho \geq 2$, we further propose a DMBA$^{++}$ strategy to directly modifies the corresponding coordinates in the density map generated by the Gaussian kernel. For example, when $\rho = 2$, we modify the values one to two at the head-related locations so that a single point now represents two head information. Unlike DBMA$^+$, DMBA$^{++}$ is performed after the generation of the density map. In this case, the density map $z'$ can be generated as:

$$z'_i (p \mid m_i) = \rho \sum_{i=1}^{c_i} \mathcal{N}^{gt} \left( p \mid \mu = P_i^j, \sigma^2 \right), \ p \in \boldsymbol{m}_i, \qquad (10)$$

**Poisoning and Inference.** DMBAs poison a small subset of the training data for both the images and their ground truth count or density map to produce a poisoned subset $\mathcal{D}_p$. The rest of the training data is kept clean, i.e., $\mathcal{D}_c$. After training on the poisoned training dataset ($\mathcal{D} = \mathcal{D}_c \cup \mathcal{D}_p$), the model will learn the correlation between the trigger pattern and the target density/count. At inference time, the attacker will attach the trigger pattern to any test image to obtain the targeted level of counting value.

## 4 EXPERIMENTS

In this section, we evaluate the effectiveness of our two DMBAs against 5 crowd counting models and conduct an ablation study of the trigger types and sizes to explain what make an effective trigger for crowd counting. The resistance to state-of-the-art backdoor defense methods are shown in later.

### 4.1 Experimental Setup

**Datasets and Networks.** We choose four benchmark counting datasets including ShanghaiTech A&B (SHA & SHB) [94], Venice [46] and TRANCOS [57] to evaluate our attacks. SHA [94] contains 482 crowd images (both RGB and Gray Scale) with a total of 241,667 annotation points, while SHB [94] contains 716 high-resolution

| Attacking Different Models on SHB Dataset | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Model →** | | | **CSRnet** | | **BayesianCC** | | **CAN** | | **SFAnet** | | **KDMG** | |
| **Attack ↓** | $\rho$ | $\gamma$ | $\hat{\rho}_{clean}$ | $\hat{\rho}_{dirty}$ | $\hat{\rho}_{clean}$ | $\hat{\rho}_{dirty}$ | $\hat{\rho}_{clean}$ | $\hat{\rho}_{dirty}$ | $\hat{\rho}_{clean}$ | $\hat{\rho}_{dirty}$ | $\hat{\rho}_{clean}$ | $\hat{\rho}_{dirty}$ |
| None | 1 | 0 | 1.00 | 1.25 | 0.98 | 0.83 | 1.03 | 0.97 | 0.99 | 0.78 | 1.02 | 0.84 |
| TriOly | 0.2 | 10% | 1.03 | 1.05 | 0.96 | 0.94 | 1.01 | 1.00 | 1.04 | 0.97 | 1.03 | 0.98 |
| DMBA⁻ | 0.2 | 5% | 1.01 | 0.33 | **0.97** | **0.23** | 0.96 | 0.38 | 1.06 | 0.44 | 1.04 | 0.39 |
| | | 10% | 0.98 | 0.24 | 0.96 | 0.16 | 0.93 | 0.29 | 1.05 | 0.59 | 1.02 | 0.29 |
| | | 15% | **0.99** | **0.22** | 0.97 | 0.16 | **0.99** | **0.28** | 1.08 | 0.41 | 1.02 | 0.28 |
| | | 20% | **1.01** | **0.22** | 0.97 | 0.14 | 0.92 | 0.28 | **1.08** | **0.34** | **1.03** | **0.26** |
| | 0.3 | 5% | 1.00 | 0.38 | 0.96 | 0.37 | 0.98 | 0.44 | **0.96** | **0.44** | 1.04 | 0.45 |
| | | 10% | 0.98 | 0.36 | 0.96 | 0.28 | 0.95 | 0.35 | 0.93 | 0.42 | 1.03 | 0.41 |
| | | 15% | 0.98 | 0.28 | **0.98** | **0.32** | **0.96** | **0.34** | 1.00 | 0.51 | 1.02 | 0.39 |
| | | 20% | **1.02** | **0.30** | 0.98 | 0.25 | 0.92 | 0.28 | 1.05 | 0.47 | **1.01** | **0.35** |
| | 0.4 | 5% | 1.02 | 0.46 | 0.98 | 0.50 | 0.93 | 0.46 | 1.06 | 0.55 | 1.01 | 0.53 |
| | | 10% | **1.02** | **0.39** | 0.99 | 0.50 | 0.94 | 0.44 | 0.94 | 0.46 | **1.04** | **0.52** |
| | | 15% | 1.06 | 0.46 | **0.98** | **0.44** | **0.95** | **0.42** | **0.99** | **0.41** | 1.05 | 0.52 |
| | | 20% | **1.00** | **0.39** | **0.97** | **0.44** | 0.96 | 0.43 | 0.93 | 0.48 | 1.00 | 0.45 |
| | 0.5 | 5% | **0.98** | **0.48** | 0.97 | 0.61 | 1.02 | 0.58 | 0.95 | 0.56 | 1.03 | 0.68 |
| | | 10% | 1.01 | 0.47 | **0.97** | **0.51** | 0.97 | 0.52 | 0.94 | 0.56 | 1.02 | 0.61 |
| | | 15% | **0.99** | **0.48** | 0.97 | 0.59 | 0.92 | 0.47 | **0.98** | **0.56** | 1.03 | 0.61 |
| | | 20% | 1.00 | 0.46 | 0.95 | 0.52 | **0.98** | **0.50** | 1.03 | 0.61 | **1.01** | **0.58** |
| DMBA⁺ | 1.2 | 10% | **1.02** | **1.13** | 0.97 | 1.10 | 0.99 | 1.18 | **1.05** | **1.14** | 0.98 | 1.17 |
| | 1.5 | 10% | **1.03** | **1.38** | 0.97 | 1.48 | 1.04 | 1.49 | 1.03 | 1.43 | 0.97 | 1.44 |
| DMBA⁺⁺ | 2 | 10% | **1.00** | **1.77** | 0.96 | 1.75 | 1.05 | 1.92 | 1.17 | 1.82 | **1.06** | **1.78** |
| | 3 | 10% | 1.00 | 2.54 | 0.97 | 2.56 | 0.99 | 3.05 | **1.23** | **3.03** | 1.05 | 2.48 |

**Table 1: Attack performance of DMBAs against 5 crowd counting models on SHB dataset. $\rho$ is the targeted manipulation ratio. Here, the poisoning rate $\gamma$ is fixed to 10% while the trigger pattern is "Rain". The most close-to-target results are boldfaced.**

crowd images captured from Shanghai streets with a total of 88,488 annotation points. Venice [46] contains 167 fixed resolution annotated frames taken from 4 different sequences, with a total of 35,440 annotation points. TRANCOS [57] contains 1244 masked images with a total of 46,734 annotated vehicles. Among the four datasets, SHA is more challenging than the other three datasets, as it has more diverse scenes and higher density. We apply our attacks on 5 counting models, including CSRnet [36], CAN [46], Bayesian Crowd Counting [54], SFA [67], and KDMG [71]. The training procedure of the backdoored models follows their original papers. Particularly, SGD [4] with learning rate 1e-7 is used to train CSRnet [36] and CAN [46], while Adam optimizer [26] with learning rate 1e-6, 1e-7, 5e-7 is used to train Bayesian Crowd couting [54], SFA [67] and KDMG [71], respectively. It is crucial to highlight that the above attack strategies are still effective even though Bayesian Crowd couting[54] follows point-supervised regression learning.

**Evaluation Metrics.** The Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE) are two standard performance metrics for crowd counting models. When applied on clean test inputs, we can obtain Clean MAE (CMAE) and Clean RMSE (CRMSE) to measure the model's performance with respect to the clean ground truth density maps. When applied on dirty (backdoored) test inputs, we can obtain Adversarial MAE (AMAE) and Adversarial RMSE (ARMSE) to measure the model's performance with respect to the altered ground truth density maps. However, MAE and RMSE measures cannot accurately reflect the closeness of the model's estimation to the targeted ratio $\rho$. To solve this problem, we further propose two new metrics as the main performance metrics for crowd counting backdoor attacks:

$$\hat{\rho}_{clean} = \frac{1}{N} \sum_{i=1}^{N} \frac{\hat{c}_i}{c_i}, \quad \hat{\rho}_{dirty} = \frac{1}{N} \sum_{i=1}^{N} \frac{\hat{c}_i}{\rho \cdot c_i}, \quad (11)$$



**Figure 4: The three types of triggers exploited in this work: "Rain", "Snow" and "Light". The triggers are blended into the poisoned image as large and defense background (highlighted by the red rectangles).**

where, $N$ is the number of test images, $\rho$ is the target ratio of the attack, and $c_i$ and $\hat{c}_i$ donate the ground truth and estimated counts, respectively. Note that $\hat{\rho}_{clean}$ is computed on clean test inputs while $\hat{\rho}_{dirty}$ is computed on backdoored test inputs. Intuitively, a $\hat{\rho}_{clean}$ close to 1 indicates good clean performance (i.e., estimated counts are close to the ground truth) while a $\hat{\rho}_{dirty}$ close to backdoor target $\rho \neq 1$ indicates effective attacks (i.e., estimated counts are close to the backdoor target). We will report the $\hat{\rho}_{clean}$ and $\hat{\rho}_{dirty}$ results in the main text and leave the CMAE/AMAE and CRMSE/ARMSE results to the appendix.

**Baselines and Attack Setup.** Since there is no existing crowd counting backdoor attacks, we take the test results of the clean-trained models on benign and backdoored test images as our baseline. For each of our attacks, we blend one of the "Rain", "Snow" or "Light" trigger pattern in the form of background into the clean training image (see Figure 4). Following Eq. (6), the pixel intensity ratio between the trigger and the clean image is set to be 3:7 so as to ensure the annotations (crowd or vehicle) are not blocked. We test multiple poisoning (injection) rates ranging from 5% to 20% on all four datasets. For our DMBA⁻ attack, we test multiple density retention rate for the SHA dataset ($\rho$ = 0.05, 0.1, 0.15 and

0.2) as well as the other 3 datasets (0.2, 0.3, 0.4 and 0.5). The varying retention rates are selected based on the annotation density of the datasets. For our DMBA+ attack, the density boosting rates $\rho = 1.2$, 1.5 are tested for the SHA dataset, while $\rho = 2, 3$ for the other three datasets. More training details and randomized ablation can be found in the supplementary material.

## 4.2 Effectiveness of our DMBAs

We first demonstrate the effectiveness of our attacks with the "Rain" trigger pattern (see Figure 4) from two perspectives: 1) the attack performance against the 5 counting models on the same dataset (i.e. SHB [57]); and 2) the attack performance against one most representative model (i.e. CSRNet) across different datasets. The $\hat{\rho}_{clean}$ and $\hat{\rho}_{dirty}$ results are reported in Table 1 and Table 2, respectively, while the CMAE/AMAE and CRMSE/ARMSE results are in Appendix A.

**Effectiveness on Different Counting Models.** One key observation from Table 1 is that, against all 5 counting models, the backdoored models by our DMBAs have achieved a $\hat{\rho}_{dirty}$ that is very close to the backdoor target (the $\rho$ column) while maintaining a $\hat{\rho}_{clean}$ that is very close to 1. This verifies the effectiveness and stealthiness of our attacks on crowd counting models, i.e., the backdoored models have been successfully controlled to output the targeted counts (whether reduced or boosted) by the trigger pattern while behaving normally on clean test data. We can also observe that stronger backdoor targets (e.g., $\rho = 0.2$ and $\rho = 2$) are more difficult to achieve than weaker targets (e.g., $\rho = 0.5$ and $\rho = 1.2$), exhibiting larger gap from the targeted $\rho$. This is because, in crowd counting, learning exceptionally sparse or dense density maps can be overrode by the learning of regular density maps, making it hard to establish the backdoor correlation with extremely small or large $\rho$. This phenomenon is different from classification backdoor attacks where the target is a hard label and the attack performance stays almost the same when choosing a different target [31, 56]. This reveals one unique challenge of attacking regression models with respect to different targeted manipulation ratios. In Table 1, it also shows the failure ($\hat{\rho}_{dirty}$ remains close to 1) of TriOly (Trigger Only) attack which only injects the trigger pattern but does not alter the ground truth labels. This confirms the importance of density map altering for attacking crowd counting models. By comparing the poisoning rates, we find that 15% poisoning is sufficient for effective attacks against the 5 models, and increasing the poisoning rate can further improve the attack but only slightly.

**Effectiveness Across Different Datasets.** Here, we apply the "Rain" trigger pattern with our DMBAs (poisoning rate $\gamma = 10\%$) to attack the CSRnet [36] on four different datasets. The results are reported in Table 2. At a high level, our attacks can consistently achieve the backdoor target (i.e., $\hat{\rho}_{dirty}$ is close to $\rho$) across the 4 datasets, with slight variations on different datasets. This is because each dataset has its own learning difficulty (indicated by the closeness of $\hat{\rho}_{clean}$ to 1) and so is the attack. The SHA [57] dataset is notably harder to attack than the other 3 datasets, showing moderately larger gap between $\hat{\rho}_{dirty}$ and $\rho$, especially when $\rho$ is extremely small ($\rho = 0.2$) or large ($\rho = 2$). It also shows a larger $\hat{\rho}_{clean}$-to-1 and $\hat{\rho}_{dirty}$-to-$\rho$ gaps on SHA than on other datasets. This indicates that datasets that consist of more diverse scenes
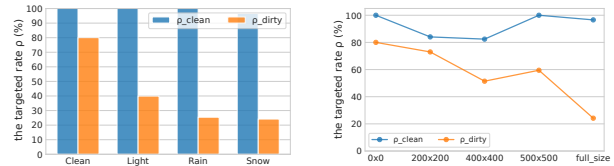


**Figure 5: Impact of Triggers. The experiments of triggers are conducted on CSRnet with SHA dataset, $\rho = 0.2$ and $\gamma = 20\%$.**

and higher densities have certain natural robustness to backdoor attacks, a phenomenon that is also different to classification backdoor attacks where high attack success rates (> 98%) can be easily achieved across different datasets [32].

By comparing the attack performance between DMBA$^-$ and DMBA$^+$/DMBA$^{++}$, we find that DMBA$^-$ is generally easier to achieve than DMBA$^+$/DMBA$^{++}$, a similar observation as in Table 1. We conjecture this is because the task of learning low-density maps is easier than learning the high-density maps where the objects are highly overlapped.

## 4.3 What Makes an Effective Trigger?

We conduct a set of ablation studies to help understand the key elements of effective crowd counting backdoor attacks. Here, we aim to gain understandings from two perspectives: 1) trigger type, and 2) trigger size. The experiments are run with CSRnet [36] on SHA dataset [57], the more challenging dataset with higher-density scenes. The backdoor target is set to $\rho = 0.2$ and the poisoning rate is fixed to $\gamma = 20\%$.

**Trigger type.** We test the three types of trigger patterns visualized in Figure 4: "Rain", "Snow" and "Light". As shown in the left subfigure of Figure 5, "Rain" and "Snow" triggers are more effective than the "Light" trigger, demonstrating closer $\hat{\rho}_{dirty}$ to 0.2 which is the attack target. This is because , compared with "Light", "Rain" and "Snow" have more scattered points that greatly interfere with the heads/objects in the original image. This experiment reveals the importance of using densely scattered trigger patterns to attack crowd counting models. Note that the trigger patterns are not restricted to natural effects, although we believe such triggers can make the attack more stealthy and are easy to simulate in real-world attacks.

**Trigger size.** Here, we take the "Snow" trigger as an example and test 5 different sizes including 0x0, 200x200, 400x400, 500x500, and the same size of the input image (see Figure 3). The results are presented in the right subfigure of Figure 5. It is evident that larger trigger patterns have a clear advantage in achieving the backdoor target $\rho = 0.2$ than small trigger patterns. This verifies the importance of using large trigger patterns to force the model to learn the backdoor correlation. Again, this is because learning the backdoor in crowd counting models is not a (relatively) independent task but rather a highly-interfered process with the original density regression task. Combining our finding in the previous experiment, we conclude that large and densely scattered background trigger patterns are the key to effective crowd counting backdoor attacks.

## 4.4 Resistance to Advanced Backdoor Defenses

Here, we test the effectiveness of three (two classic and one advanced) backdoor defense methods developed for classification models against our DMBA$^-$ attack: 1) Pruning [13], 2) Fine-Pruning

Yuhua Sun, Tailai Zhang, Xingjun Ma, Pan Zhou, Jian Lou, Zichuan Xu, Xing Di, Yu Cheng, and Lichao Sun

| Attacking Performance on Different Datasets | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Dataset → | | | SHA Dataset | | Dataset → | | | SHB Dataset | | Venice Dataset | | TRANCOS Dataset | |
| Attack ↓ | $\rho$ | $\gamma$ | $\hat{\rho}_{clean}$ | $\hat{\rho}_{dirty}$ | Attack ↓ | $\rho$ | $\gamma$ | $\hat{\rho}_{clean}$ | $\hat{\rho}_{dirty}$ | $\hat{\rho}_{clean}$ | $\hat{\rho}_{dirty}$ | $\hat{\rho}_{clean}$ | $\hat{\rho}_{dirty}$ |
| None | 1 | 0 | 1.04 | 0.80 | None | 1 | 0 | 1.00 | 1.25 | 0.91 | 0.77 | 1.00 | 1.03 |
| TriOly | 0.2 | 10% | 1.02 | 1.05 | TriOly | 0.2 | 10% | 1.03 | 1.05 | 0.93 | 0.92 | 0.99 | 0.97 |
| DMBA⁻ | 0.05 | 5% | 1.00 | 0.19 | DMBA⁻ | 0.2 | 5% | 1.01 | 0.33 | 0.85 | 0.34 | 0.95 | 0.45 |
| | | 10% | 1.01 | 0.14 | | | 10% | 0.98 | 0.24 | 0.84 | 0.25 | 0.94 | 0.39 |
| | | 15% | 1.01 | 0.11 | | | 15% | **0.99** | **0.22** | **0.84** | **0.21** | **0.92** | **0.33** |
| | | 20% | **1.01** | **0.10** | | | 20% | 1.01 | 0.22 | 0.83 | 0.23 | 0.94 | 0.39 |
| | 0.1 | 5% | 1.02 | 0.24 | | 0.3 | 5% | 1.00 | 0.38 | 0.90 | 0.44 | 0.97 | 0.59 |
| | | 10% | 1.02 | 0.19 | | | 10% | 0.98 | 0.36 | 0.87 | 0.33 | 0.90 | 0.48 |
| | | 15% | **1.02** | **0.18** | | | 15% | 0.98 | 0.28 | 0.88 | 0.31 | 0.92 | 0.42 |
| | | 20% | 1.02 | 0.19 | | | 20% | 1.02 | 0.30 | 0.87 | 0.30 | 0.94 | 0.36 |
| | 0.15 | 5% | 1.02 | 0.34 | | 0.4 | 5% | 1.02 | 0.46 | 0.88 | 0.43 | 0.94 | 0.65 |
| | | 10% | 1.00 | 0.29 | | | 10% | **1.02** | **0.39** | 0.84 | 0.37 | 0.94 | 0.53 |
| | | 15% | **1.02** | **0.24** | | | 15% | 1.06 | 0.46 | **0.84** | **0.39** | 0.93 | 0.49 |
| | | 20% | **1.01** | **0.24** | | | 20% | **1.00** | **0.39** | 0.85 | 0.37 | **0.92** | **0.48** |
| | 0.2 | 5% | 1.02 | 0.33 | | 0.5 | 5% | **0.98** | **0.48** | **0.90** | **0.53** | 0.95 | 0.74 |
| | | 10% | 1.03 | 0.29 | | | 10% | 1.01 | 0.47 | 0.85 | 0.44 | 0.94 | 0.63 |
| | | 15% | 1.01 | 0.26 | | | 15% | **0.99** | **0.48** | 0.88 | 0.44 | 0.93 | 0.62 |
| | | 20% | **1.01** | **0.25** | | | 20% | 1.00 | 0.46 | 0.87 | 0.46 | **1.02** | **0.56** |
| DMBA⁺ | 1.2 | 10% | **1.04** | **1.12** | DMBA⁺ | 1.2 | 10% | **1.02** | **1.13** | 0.98 | 1.11 | **1.00** | **1.15** |
| | 1.5 | 10% | 1.07 | 1.38 | | 1.5 | 10% | 1.03 | 1.38 | 1.02 | 1.48 | 0.96 | 1.35 |
| DMBA⁺⁺ | 2 | 10% | **1.08** | **1.84** | DMBA⁺⁺ | 2 | 10% | 1.00 | 1.77 | 1.01 | 1.72 | 0.97 | 2.01 |
| | 3 | 10% | 1.10 | 2.61 | | 3 | 10% | 1.00 | 2.54 | 1.05 | 2.35 | 0.98 | 2.72 |

**Table 2: Attack performance of DMBAs against CSRnet on different datasets. $\rho$ is the targeted manipulation ratio. Here, the poisoning rate $\gamma$ is fixed to 10% while the trigger pattern is "Rain". The most close-to-target results are boldfaced.**

[42], and 3) Adversarial Neural Pruning (ANP) [82]. We run these experiments with CSRnet model on SHA dataset with trigger pattern "Rain", backdoor target $\rho = 0.2$ and poisoning rate $\gamma = 10\%$.

**Pruning and Fine-Pruning.** Both methods prune the neurons that stay dormant on clean inputs but can potentially be triggered by backdoored inputs to mitigate backdoor attacks. Here, we apply pruning and fine-pruning on the DMBA⁻-backdoored CSRnet to prune neurons from the last two layers of both the back-end and front-end modules of the network. We follow the *prune-then-test* pipeline for pruning and the *prune-finetune-then-test* pipeline for Fine-pruning. We prune the the selected layers until 90% of the neurons are removed. We find that both defenses can mitigate our DMBA⁻ attack to certain extent, yet they both significantly degrade the model's performance on the clean inputs. Particularly, when $\hat{\rho}_{dirty}$ was recovered to 0.95 (originally close to the backdoor target 0.2), the $\hat{\rho}_{dirty}$ of the model increases to 1.83 (originally close to 1) while the MAE increases from 92.41 to 415.55. More detailed results can be found in Appendix B.

**Adversarial Neuron Pruning (ANP).** ANP is one of the state-of-the-art defense methods that locates and prunes backdoor neurons based on the neurons' sensitivity to adversarial perturbations [82]. ANP was originally proposed for classification models. Here we adapt ANP to defend crowd counting models. We first optimized the neural perturbation as follows:

$$\max_{\delta, \xi \in [-\epsilon, \epsilon]^n} \mathcal{L}_{\mathcal{D}_v} = \frac{1}{2N_v} \left\| \hat{z}_i - z_i^{gt} \right\|_2^2, \quad (12)$$

where, $\delta$, $\xi$ and $\epsilon$ are the 3 hyper-parameters of ANP [82], $n$ is the number of neurons, and $N_v$ is the number of clean samples from a small clean validation set $\mathcal{D}_v$. ANP then locates backdoor neurons by learning a mask via the following:

$$\min_{\mathbf{m} \in [0,1]^n} \left[ \alpha \mathcal{L}_{\mathcal{D}_v}(\mathbf{m} \odot \theta) + (1-\alpha) \cdot \max_{\delta, \xi} \mathcal{L}_{\mathcal{D}_v}([\delta, \xi] \odot \theta) \right], \quad (13)$$

where, $\theta$ are the model parameters, $\mathbf{m}$ is the neuron mask and $\alpha$ is a trade-off coefficient [82]. After 2000 iterations of pruning following the setting in [82], $\hat{\rho}_{dirty}$ of the backdoored CSRnet model is still 0.25, which is close to no ANP defense ($\hat{\rho}_{dirty} = 0.21$). Moreover, the results on the clean test data show that ANP greatly degrade the model's clean performance (see Figure **??** in Appendix B).

The above results indicate that backdoored neurons are mixed with clean neurons in the backdoored models by our DMBA⁻ attack, making it hard to segregate and prune backdoor neurons without hurting the clean performance. This means that our DMBA⁻ attack is fairly resistant to these defense methods.

## 5 CONCLUSION

In this paper, we studied the problem of backdoor attack on crowd counting models. We first verified the ineffectiveness of classification backdoor attacks on crowd counting models, then proposed two novel Density Manipulation Backdoor Attacks (DMBA⁻ and DMBA⁺) to attack crowd counting models to produce targeted count estimations. We demonstrated the effectiveness of our DMBA attacks on 5 crowd counting models and 4 datasets with low poisoning rate. We also provide an analysis of the key elements of effective attacks: 1) large and dense trigger patterns, and 2) the alteration of the ground truth density maps. We hope our attacks can serve as strong baselines for the robustness evaluation of crowd counting models and development of effective defenses [9, 98].

## 6 ACKNOWLEDGMENTS

# REFERENCES

[1] Prithvi N. Amin, Sayali S. Moghe, Sparsh N. Prabhakar, and Charusheela Nehete. 2021. Deep Learning Based Face Mask Detection and Crowd Counting. *I2CT* (2021), 1–5.

[2] Mauro Barni, Kassem Kallas, and Benedetta Tondi. 2019. A New Backdoor Attack in CNNS by Training Set Corruption Without Label Poisoning. *2019 IEEE International Conference on Image Processing (ICIP)* (2019), 101–105.

[3] Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao. 2020. Yolov4: Optimal speed and accuracy of object detection. *arXiv preprint arXiv:2004.10934* (2020).

[4] Léon Bottou. 2010. Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT'2010*. Springer, 177–186.

[5] Xinkun Cao, Zhipeng Wang, Yanyun Zhao, and Fei Su. 2018. Scale aggregation network for accurate and efficient crowd counting. In *ECCV*. 734–750.

[6] Antoni B. Chan and Nuno Vasconcelos. 2009. Bayesian Poisson regression for crowd counting. *ICCV* (2009), 545–551.

[7] Bryant Chen, Wilka Carvalho, Nathalie Baracaldo, Heiko Ludwig, Benjamin Edwards, Taesung Lee, Ian Molloy, and Biplav Srivastava. 2018. Detecting backdoor attacks on deep neural networks by activation clustering. *arXiv preprint arXiv:1811.03728* (2018).

[8] Jinyin Chen, Haibin Zheng, Mengmeng Su, Tianyu Du, Changting Lin, and Shouling Ji. 2019. Invisible Poisoning: Highly Stealthy Targeted Poisoning Attack. In *ICISC*.

[9] Tianlong Chen, Sijia Liu, Shiyu Chang, Yu Cheng, Lisa Amini, and Zhangyang Wang. 2020. Adversarial robustness: From self-supervised pre-training to fine-tuning. In *CVPR*. 699–708.

[10] Xinyun Chen, Chang Liu, Bo Li, Kimberly Lu, and Dawn Song. 2017. Targeted backdoor attacks on deep learning systems using data poisoning. *arXiv preprint arXiv:1712.05526* (2017).

[11] Siyuan Cheng, Yingqi Liu, Shiqing Ma, and Xiangyu Zhang. 2021. Deep Feature Space Trojan Attack of Neural Networks by Controlled Detoxification. In *AAAI*.

[12] Ronan Collobert, Koray Kavukcuoglu, and Clément Farabet. 2011. Torch7: A Matlab-like Environment for Machine Learning. In *NIPS*.

[13] Guneet Singh Dhillon, Kamyar Azizzadenesheli, Zachary Chase Lipton, Jeremy Bernstein, Jean Kossaifi, Aran Khanna, and Anima Anandkumar. 2018. Stochastic Activation Pruning for Robust Adversarial Defense. *ArXiv* abs/1803.01442 (2018).

[14] Kevin Eykholt, I. Evtimov, Earlence Fernandes, Bo Li, Amir Rahmati, Chaowei Xiao, Atul Prakash, Tadayoshi Kohno, and Dawn Xiaodong Song. 2018. Robust Physical-World Attacks on Deep Learning Visual Classification. *CVPR* (2018), 1625–1634.

[15] Pedro F Felzenszwalb, Ross B Girshick, David McAllester, and Deva Ramanan. 2009. Object detection with discriminatively trained part-based models. *IEEE transactions on pattern analysis and machine intelligence* 32, 9 (2009), 1627–1645.

[16] Yansong Gao, Change Xu, Derui Wang, Shiping Chen, Damith C Ranasinghe, and Surya Nepal. 2019. Strip: A defence against trojan attacks on deep neural networks. In *ACSAC*.

[17] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. 2015. Explaining and Harnessing Adversarial Examples. *ICLR* abs/1412.6572 (2015).

[18] Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. 2017. BadNets: Identifying Vulnerabilities in the Machine Learning Model Supply Chain. *ArXiv* abs/1708.06733 (2017).

[19] Ricardo Guerrero-Gómez-Olmedo, Beatriz Torre-Jiménez, Roberto Javier López-Sastre, Saturnino Maldonado-Bascón, and Daniel Oñoro-Rubio. 2015. Extremely Overlapping Vehicle Counting. In *IbPRIA*.

[20] Shamim Hossain, Hamid A. Jalab, Fariha Zulfiqar, and Mahfuza Pervin. 2019. Renal Cancer Cell Nuclei Detection from Cytological Images Using Convolutional Neural Network for Estimating Proliferation Rate. *Journal of Telecommunication, Electronic and Computer Engineering* 11 (2019), 63–71.

[21] Hongsheng Hu, Zoran A. Salcic, Gillian Dobbie, Jinjun Chen, Lichao Sun, and Xuyun Zhang. 2022. Membership Inference via Backdooring. *ArXiv* abs/2206.04823 (2022).

[22] Haroon Idrees, Imran Saleemi, Cody Seibert, and Mubarak Shah. 2013. Multi-source multi-scale counting in extremely dense crowd images. In *CVPR*. 2547–2554.

[23] Haroon Idrees, Muhmmad Tayyab, Kishan Athrey, Dong Zhang, Somaya Ali Al-Maadeed, Nasir M. Rajpoot, and Mubarak Shah. 2018. Composition Loss for Counting, Density Map Estimation and Localization in Dense Crowds. *ArXiv* abs/1808.01050 (2018).

[24] Xiaoheng Jiang, Li Zhang, Mingliang Xu, Tianzhu Zhang, Pei Lv, Bing Zhou, Xin Yang, and Yanwei Pang. 2020. Attention scaling for crowd counting. In *CVPR*. 4706–4715.

[25] Di Kang, Debarun Dhar, and Antoni B. Chan. 2020. Incorporating Side Information by Adaptive Convolution. *International Journal of Computer Vision* (2020), 1–22.

[26] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).

[27] Issam H. Laradji, Negar Rostamzadeh, Pedro H. O. Pinheiro, David Vázquez, and Mark W. Schmidt. 2018. Where are the Blobs: Counting by Localization with Point Supervision. *ArXiv* abs/1807.09856 (2018).

[28] Victor S. Lempitsky and Andrew Zisserman. 2010. Learning To Count Objects in Images. In *NIPS*.

[29] Shaofeng Li, Minhui Xue, Benjamin Zhao, Haojin Zhu, and Xinpeng Zhang. 2020. Invisible backdoor attacks on deep neural networks via steganography and regularization. *IEEE Transactions on Dependable and Secure Computing* (2020).

[30] Shaofeng Li, Minhui Xue, Benjamin Zi Hao Zhao, Haojin Zhu, and Xinpeng Zhang. 2019. Invisible Backdoor Attacks on Deep Neural Networks via Steganography and Regularization. *arXiv preprint arXiv:1909.02742* (2019).

[31] Yuezun Li, Yiming Li, Baoyuan Wu, Longkang Li, Ran He, and Siwei Lyu. 2021. Invisible Backdoor Attack With Sample-Specific Triggers. In *ICCV*. 16463–16472.

[32] Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, and Xingjun Ma. 2021. Anti-Backdoor Learning: Training Clean Models on Poisoned Data. *NeurIPS* (2021).

[33] Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, and Xingjun Ma. 2021. Neural attention distillation: Erasing backdoor triggers from deep neural networks. *ICLR* (2021).

[34] Yiming Li, Baoyuan Wu, Yong Jiang, Zhifeng Li, and Shu-Tao Xia. 2020. Backdoor learning: A survey. *arXiv preprint arXiv:2007.08745* (2020).

[35] Yiming Li, Tongqing Zhai, Yong Jiang, Zhifeng Li, and Shu-Tao Xia. 2021. Backdoor Attack in the Physical World. *arXiv preprint arXiv:2104.02361* (2021).

[36] Yuhong Li, Xiaofan Zhang, and Deming Chen. 2018. CSRNet: Dilated Convolutional Neural Networks for Understanding the Highly Congested Scenes. *CVPR* (2018), 1091–1100.

[37] Yiming Li, Haoxiang Zhong, Xingjun Ma, Yong Jiang, and Shu-Tao Xia. 2022. Few-shot backdoor attacks on visual object tracking. *ICLR* (2022).

[38] Dongze Lian, Jing Li, Jia Zheng, Weixin Luo, and Shenghua Gao. 2019. Density map regression guided detection network for rgb-d crowd counting and localization. In *CVPR*. 1821–1830.

[39] Cong Liao, Haoti Zhong, Anna Squicciarini, Sencun Zhu, and David Miller. 2020. Backdoor embedding in convolutional neural network models via invisible perturbation. *CODASPY* (2020).

[40] Zhe L. Lin and Larry S. Davis. 2010. Shape-Based Human Detection and Segmentation via Hierarchical Part-Template Matching. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 32 (2010), 604–618.

[41] Jiang Liu, Chenqiang Gao, Deyu Meng, and Alexander Hauptmann. 2018. DecideNet: Counting Varying Density Crowds Through Attention Guided Detection and Density Estimation. *CVPR* (2018), 5197–5206.

[42] Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. 2018. Fine-pruning: Defending against backdooring attacks on deep neural networks. In *International Symposium on Research in Attacks, Intrusions, and Defenses*. Springer, 273–294.

[43] Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. 2018. Fine-pruning: Defending against backdooring attacks on deep neural networks. In *RAID*. Springer, 273–294.

[44] Lingbo Liu, Jiaqi Chen, Hefeng Wu, Guanbin Li, Chenglong Li, and Liang Lin. 2021. Cross-Modal Collaborative Representation Learning and a Large-Scale RGBT Benchmark for Crowd Counting. In *CVPR*. 4823–4833.

[45] Lingbo Liu, Zhilin Qiu, Guanbin Li, Shufan Liu, Wanli Ouyang, and Liang Lin. 2019. Crowd counting with deep structured scale integration network. In *ICCV*. 1774–1783.

[46] Weizhe Liu, Mathieu Salzmann, and Pascal Fua. 2019. Context-Aware Crowd Counting. *CVPR* (2019), 5094–5103.

[47] Weizhe Liu, Mathieu Salzmann, and P. Fua. 2019. Using Depth for Pixel-Wise Detection of Adversarial Attacks in Crowd Counting. *ArXiv* abs/1911.11484 (2019).

[48] Xinyan Liu, Guorong Li, Zhenjun Han, Weigang Zhang, Yifan Yang, Qingming Huang, and Nicu Sebe. 2021. Exploiting sample correlation for crowd counting with multi-expert network. In *CVPR*. 3215–3224.

[49] Xialei Liu, Joost Van De Weijer, and Andrew D Bagdanov. 2018. Leveraging unlabeled data for crowd counting by learning to rank. In *CVPR*. 7661–7669.

[50] Yingqi Liu, Shiqing Ma, Yousra Aafer, Wen-Chuan Lee, Juan Zhai, Weihang Wang, and Xiangyu Zhang. 2017. Trojaning attack on neural networks. (2017).

[51] Yunfei Liu, Xingjun Ma, James Bailey, and Feng Lu. 2020. Reflection backdoor: A natural backdoor attack on deep neural networks. In *ECCV*. Springer, 182–199.

[52] Xingjun Ma, Bo Li, Yisen Wang, Sarah Monazam Erfani, Sudanthi N. R. Wijewickrema, Michael E. Houle, Grant Robert Schoenebeck, Dawn Xiaodong Song, and James Bailey. 2018. Characterizing Adversarial Subspaces Using Local Intrinsic Dimensionality. *ICLR* (2018).

[53] Zhiheng Ma, Xiaopeng Hong, Xing Wei, Yunfeng Qiu, and Yihong Gong. 2021. Towards a universal model for cross-dataset crowd counting. In *CVPR*. 3205–3214.

[54] Zhiheng Ma, Xing Wei, Xiaopeng Hong, and Yihong Gong. 2019. Bayesian loss for crowd count estimation with point supervision. In *ICCV*. 6142–6151.

[55] Yanda Meng, Hongrun Zhang, Yitian Zhao, Xiaoyun Yang, Xuesheng Qian, Xiaowei Huang, and Yalin Zheng. 2021. Spatial Uncertainty-Aware Semi-Supervised Crowd Counting. In *CVPR*. 15549–15559.

[56] Anh Nguyen and Anh Tran. 2020. Input-aware dynamic backdoor attack. In *NeurIPS*.

[57] Daniel Oñoro-Rubio and Roberto Javier López-Sastre. 2016. Towards Perspective-Free Object Counting with Deep Learning. In *ECCV*.

[58] Viresh Ranjan, Hieu Le, and Minh Hoai. 2018. Iterative crowd counting. In *ECCV*. 270–285.

[59] Aniruddha Saha, Akshayvarun Subramanya, and Hamed Pirsiavash. 2020. Hidden trigger backdoor attacks. In *AAAI*, Vol. 34. 11957–11965.

[60] Ahmed Salem, Michael Backes, and Yang Zhang. 2020. Don't Trigger Me! A Triggerless Backdoor Attack Against Deep Neural Networks. *arXiv preprint arXiv:2010.03282* (2020).

[61] Deepak Babu Sam, Shiv Surya, and R. Venkatesh Babu. 2017. Switching Convolutional Neural Network for Crowd Counting. *CVPR* (2017), 4031–4039.

[62] Ali Shafahi, W Ronny Huang, Mahyar Najibi, Octavian Suciu, Christoph Studer, Tudor Dumitras, and Tom Goldstein. 2018. Poison frogs! targeted clean-label poisoning attacks on neural networks. In *NeurIPS*.

[63] Zan Shen, Yi Xu, Bingbing Ni, Minsi Wang, Jianguo Hu, and Xiaokang Yang. 2018. Crowd counting via adversarial cross-scale consistency pursuit. In *CVPR*. 5245–5254.

[64] Vishwanath A. Sindagi and Vishal M. Patel. 2017. CNN-Based cascaded multi-task learning of high-level prior and density estimation for crowd counting. *AVSS* (2017), 1–6.

[65] Jacob Steinhardt, Pang Wei Koh, and Percy Liang. 2017. Certified Defenses for Data Poisoning Attacks. In *NIPS*.

[66] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, D. Erhan, Ian J. Goodfellow, and Rob Fergus. 2014. Intriguing properties of neural networks. *ICLR* (2014).

[67] Pongpisit Thanasutives, Ken ichi Fukui, Masayuki Numao, and Boonserm Kijsirikul. 2021. Encoder-Decoder Based Convolutional Neural Networks with Multi-Scale-Aware Modules for Crowd Counting. *ICPR* (2021), 2382–2389.

[68] Brandon Tran, Jerry Li, and Aleksander Madry. 2018. Spectral signatures in backdoor attacks. In *NeurIPS*.

[69] Alexander Turner, Dimitris Tsipras, and Aleksander Madry. 2019. Label-consistent backdoor attacks. *arXiv preprint arXiv:1912.02771* (2019).

[70] Immanuel Jose C Valencia, Elmer P Dadios, Alexis M Fillone, John Carlo V Puno, Renann G Baldovino, and Robert Kerwin C Billones. 2021. Vision-based Crowd Counting and Social Distancing Monitoring using Tiny-YOLOv4 and DeepSORT. In *ISC2*. IEEE, 1–7.

[71] Jia Wan, Qingzhong Wang, and Antoni B. Chan. 2022. Kernel-Based Density Map Generation for Dense Object Counting. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 44 (2022), 1357–1370.

[72] Boyu Wang, Huidong Liu, Dimitris Samaras, and Minh Hoai Nguyen. 2020. Distribution matching for crowd counting. *Advances in Neural Information Processing Systems* 33 (2020), 1595–1607.

[73] Qi Wang, Junyu Gao, Wei Lin, and Yuan Yuan. 2019. Learning from synthetic data for crowd counting in the wild. In *ICCV*. 8198–8207.

[74] Xin Wang, Bin Wang, and Liming Zhang. 2011. Airport Detection in Remote Sensing Images Based on Visual Attention. In *ICONIP*.

[75] Yisen Wang, Xingjun Ma, James Bailey, Jinfeng Yi, Bowen Zhou, and Quanquan Gu. 2019. On the Convergence and Robustness of Adversarial Training. In *ICML*.

[76] Yisen Wang, Difan Zou, Jinfeng Yi, James Bailey, Xingjun Ma, and Quanquan Gu. 2020. Improving Adversarial Robustness Requires Revisiting Misclassified Examples. In *ICLR*.

[77] Longyin Wen, Dawei Du, Pengfei Zhu, Qinghua Hu, Qilong Wang, Liefeng Bo, and Siwei Lyu. 2021. Detection, tracking, and counting meets drones in crowds: A benchmark. In *CVPR*. 7812–7821.

[78] Emily Wenger, Josephine Passananti, Arjun Nitin Bhagoji, Yuanshun Yao, Haitao Zheng, and Ben Y Zhao. 2021. Backdoor Attacks Against Deep Learning Systems in the Physical World. In *CVPR*. 6206–6215.

[79] Emily Wenger, Josephine Passananti, Yuanshun Yao, Haitao Zheng, and Ben Y Zhao. 2020. Backdoor attacks on facial recognition in the physical world. *arXiv e-prints* (2020), arXiv–2006.

[80] Bo Wu and Ramakant Nevatia. 2005. Detection of multiple, partially occluded humans in a single image by Bayesian combination of edgelet part detectors. *ICCV* 1 (2005), 90–97 Vol. 1.

[81] Bo Wu and Ram Nevatia. 2007. Detection and tracking of multiple, partially occluded humans by bayesian combination of edgelet based part detectors. *International Journal of Computer Vision* 75, 2 (2007), 247–266.

[82] Dongxian Wu and Yisen Wang. 2021. Adversarial neuron pruning purifies backdoored deep models. *Advances in Neural Information Processing Systems* 34 (2021).

[83] Dongxian Wu, Yisen Wang, Shutao Xia, James Bailey, and Xingjun Ma. 2020. Skip Connections Matter: On the Transferability of Adversarial Examples Generated with ResNets. *ICLR* (2020).

[84] Qiangqiang Wu, Jia Wan, and Antoni B Chan. 2021. Dynamic Momentum Adaptation for Zero-Shot Cross-Domain Crowd Counting. In *ACM MM*. 658–666.

[85] Qiming Wu, Zhikang Zou, Pan Zhou, Xiaoqing Ye, Binghui Wang, and Ang Li. 2021. Towards Adversarial Patch Analysis and Certified Defense against Crowd Counting. *ACM MM* (2021).

[86] Chulin Xie, Keli Huang, Pin-Yu Chen, and Bo Li. 2020. Dba: Distributed backdoor attacks against federated learning. In *ICLR*.

[87] Haipeng Xiong, Hao Lu, Chengxin Liu, Liang Liu, Zhiguo Cao, and Chunhua Shen. 2019. From Open Set to Closed Set: Counting Objects by Spatial Divide-and-Conquer. *ICCV* (2019), 8361–8370.

[88] Yanyu Xu, Ziming Zhong, Dongze Lian, Jing Li, Zhengxin Li, Xinxing Xu, and Shenghua Gao. 2021. Crowd Counting With Partial Annotations in an Image. In *CVPR*. 15570–15579.

[89] Mingfu Xue, Can He, Shichang Sun, Jian Wang, and Weiqiang Liu. 2021. Robust Backdoor Attacks against Deep Neural Networks in Real Physical World. *arXiv preprint arXiv:2104.07395* (2021).

[90] Yuanshun Yao, Huiying Li, Haitao Zheng, and Ben Y Zhao. 2019. Latent backdoor attacks on deep neural networks. In *SIGSAC*. 2041–2055.

[91] Cong Zhang, Hongsheng Li, Xiaogang Wang, and Xiaokang Yang. 2015. Cross-scene crowd counting via deep convolutional neural networks. *CVPR* (2015), 833–841.

[92] Qi Zhang, Wei Lin, and Antoni B Chan. 2021. Cross-View Cross-Scene Multi-View Crowd Counting. In *CVPR*. 557–567.

[93] Shanghang Zhang, Guanhang Wu, João Paulo Costeira, and José M. F. Moura. 2017. FCN-rLSTM: Deep Spatio-Temporal Neural Networks for Vehicle Counting in City Cameras. *ICCV* (2017), 3687–3696.

[94] Yingying Zhang, Desen Zhou, Siqin Chen, Shenghua Gao, and Yi Ma. 2016. Single-Image Crowd Counting via Multi-Column Convolutional Neural Network. *CVPR* (2016), 589–597.

[95] Zhiyuan Zhang, Lingjuan Lyu, Weiqiang Wang, Lichao Sun, and Xu Sun. 2021. How to Inject Backdoors with Better Consistency: Logit Anchoring on Clean Data. *ArXiv* abs/2109.01300 (2021).

[96] Zhaoxiang Zhang, Mo Wang, and Xin Geng. 2015. Crowd counting in public video surveillance by label distribution learning. *Neurocomputing* 166 (2015), 151–163.

[97] Shihao Zhao, Xingjun Ma, Xiang Zheng, James Bailey, Jingjing Chen, and Yu-Gang Jiang. 2020. Clean-label backdoor attacks on video recognition models. In *CVPR*. 14443–14452.

[98] Chen Zhu, Yu Cheng, Zhe Gan, Siqi Sun, Tom Goldstein, and Jingjing Liu. 2020. FreeLB: Enhanced Adversarial Training for Natural Language Understanding. In *ICLR*.

[99] Chen Zhu, W Ronny Huang, Hengduo Li, Gavin Taylor, Christoph Studer, and Tom Goldstein. 2019. Transferable clean-label poisoning attacks on deep neural nets. In *ICML*. PMLR, 7614–7623.

[100] Liang Zhu, Zhijian Zhao, Chao Lu, Yining Lin, Yao Peng, and Tangren Yao. 2019. Dual Path Multi-Scale Fusion Networks with Attention for Crowd Counting. *ArXiv* abs/1902.01115 (2019).

[101] Zhikang Zou, Yu Cheng, Xiaoye Qu, Shouling Ji, Xiaoxiao Guo, and Pan Zhou. 2019. Attend to count: Crowd counting with adaptive capacity multi-scale CNNs. *Neurocomputing* 367 (2019), 75–83.

[102] Zhikang Zou, Xiaoye Qu, Pan Zhou, Shuangjie Xu, Xiaoqing Ye, Wenhao Wu, and Jin Ye. 2021. Coarse to fine: Domain adaptive crowd counting via adversarial scoring network. In *ACM MM*. 2185–2194.