# Model-Guardian: Protecting against Data-Free Model Stealing Using Gradient Representations and Deceptive Predictions*

Yunfei Yang[1,2,3], Xiaojun Chen[1,2,3,†], Yuexin Xuan[1,2,3], and Zhendong Zhao[1,2]

[1]Institute of Information Engineering, Chinese Academy of Sciences, Beijing, China
[2]State Key Laboratory of Cyberspace Security Defense, Beijing, China
[3]School of Cyber Security, University of Chinese Academy of Sciences, Beijing, China
{yangyunfei, chenxiaojun, xuanyuexin, zhaozhendong}@iie.ac.cn

*Abstract*—**Model stealing attack is increasingly threatening the confidentiality of machine learning models deployed in the cloud. Recent studies reveal that adversaries can exploit data synthesis techniques to steal machine learning models even in scenarios devoid of real data, leading to data-free model stealing attacks. Existing defenses against such attacks suffer from limitations, including poor effectiveness, insufficient generalization ability, and low comprehensiveness. In response, this paper introduces a novel defense framework named Model-Guardian. Comprising two components, Data-Free Model Stealing Detector (DFMS-Detector) and Deceptive Predictions (DPreds), Model-Guardian is designed to address the shortcomings of current defenses with the help of the artifact properties of synthetic samples and gradient representations of samples. Extensive experiments on seven prevalent data-free model stealing attacks showcase the effectiveness and superior generalization ability of Model-Guardian, outperforming eleven defense methods and establishing a new state-of-the-art performance. Notably, this work pioneers the utilization of various GANs and diffusion models for generating highly realistic query samples in attacks, with Model-Guardian demonstrating accurate detection capabilities.**

*Index Terms*—**AI Security, Model Stealing, Defense**

## I. INTRODUCTION

With the rapid growth of deep learning, some companies offer Machine Learning as a Service (MLaaS), deploying pretrained models in the cloud for commercial profit. However, vulnerabilities in MLaaS have been exposed, particularly regarding model stealing attacks [1]–[3]. In these attacks, adversaries query APIs with crafted samples to acquire annotated data for training equivalent models [4]. Attackers can then use the stolen model for adversarial attacks or membership inference attacks. Recent studies show that even without real data or knowledge of the original distribution for training data, attackers can still achieve high attack performance using data-free stealing methods [5]–[8], a trend that has gained significant attention in the research community.

To mitigate the risk of model stealing attacks, various defenses have been proposed, categorized as active and passive. Active defenses [9]–[11] perturb outputs during the prediction phase to hinder attackers from replicating the victim model's functionality. Passive defenses [12], [13] detect anomalies in query sequences, leading to service termination. However, current defenses have limitations against data-free stealing attacks: (1) Most lack data-free attack considerations, as shown in our experiments, making them ineffective and limited in generalization. (2) Active defenses cannot distinguish and terminate malicious queries, allowing attackers to generate clone models with acceptable performance by inflating queries, while reducing benign query accuracy. (3) Passive defenses, which rely on query distribution characteristics, fail to detect individual samples and could improve in detection accuracy and false positive rates (FPR).

To address these limitations, we propose Model-Guardian, a comprehensive defense safeguarding model privacy from data-free model stealing attacks. It consists of Data-Free Model Stealing Detector (DFMS-Detector) and Deceptive Predictions (DPreds) modules. Due to the varied generative model structures and diverse generated samples in data-free model stealing, we adopt more generalized representations, gradients, as training data for DFMS-Detector. In the DPreds module, we strategically perturb model outputs to maintain class probability relationships while providing users with perturbed results. During implementation, Model-Guardian adapts its responses based on query detection results. The proposed Model-Guardian exhibits several advantages: (1) Effectiveness: It can effectively defend against prevalent data-free stealing attacks. (2) Generalization: It performs well against novel attacks and data from various generative models. (3) Harmlessness: It minimally impacts benign users, preserving model accuracy.
**Contributions.** Our contributions encompass four aspects:
- Proposal of Model-Guardian, an innovative and effective defense addressing data-free model stealing attacks. This work represents a pioneering effort, specifically targeting data-free scenarios and considering adversaries' use of GANs and diffusion models for realistic data synthesis.
- Enhancement of generalization through the use of gra-

dient representations in training an ensemble detector, DFMS-Detector, within Model-Guardian. This detector accurately discerns malicious queries from multiple attacks, demonstrating low FPR on benign samples.

- Introduction of a simple yet effective prediction perturbation algorithm, Deceptive Predictions (DPreds), capable of disrupting the probability distribution and adversary's clone training process without altering the magnitude relationship of original class probabilities. Importantly, it does not impact the use of misjudged benign users.
- Conduct of extensive experiments across multiple datasets to evaluate the effectiveness of Model-Guardian. Results showcase superior performance compared to baselines in defending against data-free stealing attacks.

## II. RELATED WORK

### A. Model Stealing Attacks

**Adversarial Generation-based Attacks.** The attacker, possessing limited training data, increases query data by synthesizing adversarial samples. For instance, JBDA [14] introduces a data augmentation technique based on Jacobian matrix.

**Proxy Data-based Attacks.** Attackers are constrained to utilizing public datasets as query samples, with an emphasis on optimizing sampling strategies to enhance stealing effectiveness. KnockoffNet [15] leverages reinforcement learning to construct a transfer set. MExMI [16] represents progress by concurrently incorporating model stealing and membership inference attacks, achieving mutual performance enhancement.

**Data-Free Attacks.** The attacker, lacking access to real data, relies on synthesizing samples through a generator. DaST [17] employs a multi-branch architecture and label-controlled loss for generator to synthesize data. MAZE [5] and DFME [1] utilize zero-order gradient estimation to compute the victim's gradient, facilitating training update of generator. DFMS [4] conducts data-free stealing attacks under hard-label settings. DisGUIDE [2] and Dual Student [6] are pioneering in using the dual clone structures for stealing. QUDA [7] introduces a novel attack leveraging a frozen GAN pre-trained with publicly irrelevant data to provide weak image priors.

Data-free model stealing attacks are more practical in real-world scenarios due to challenges in obtaining victim training data and the gap between public datasets and victim's private training sets. Therefore, our research focuses on developing defenses against data-free model stealing.

### B. Defenses against Model Stealing Attacks

**Active Defenses.** Defenders add perturbations or introduce randomness to model outputs, diminishing the accuracy of clone training. RS [18] introduces fuzziness through shrewd noise addition to output probabilities. MAD [19] disrupts adversary-obtained gradients by adding perturbations to the original output. AM [9] detects OOD input presence, returning a misinformative prediction to adversary. EDM [20] introduces randomness into the model output via an ensemble of diverse models. NT [21] utilizes a Nasty Teacher to prevent model distillation and stealing. PoW [22] significantly increases adversary's query costs based on proof of work. APGP [10]

ensures complete accuracy preservation for black-box model privacy. InI [11] directly trains defensive models by isolating the adversary's training gradient.

**Passive Defenses.** Defenders discern model stealing attacks by identifying abnormal behavior in continuous query samples, followed by implementing subsequent defense measures. PRADA [12] detects malicious queries by assessing the statistical distribution of the minimum $l_2$-norm distance between queries. SEAT [23] employs a Similarity Encoder to encode visually similar samples with similar inter-code distances, determining malicious behavior based on the pair number of similar samples. FDINET [13] introduces Feature Distortion Index (FDI), enabling defenders to train a binary detector for detecting the presence of stealing attacks using FDI.

Existing active defenses struggle with impacting benign queries and only reduce clone model performance without preventing attacks. Passive defenses suffer from low detection accuracy, high false positives, and susceptibility to bypassing via discontinuous malicious queries. Moreover, there is no dedicated defense for data-free model stealing, which is the primary focus of our work.

## III. DEFENSE STRATEGY: MODEL-GUARDIAN

### A. Threat Model

**Adversary's Goals and Capabilities.** Our work focuses on model functionality stealing in image classification, wherein adversary aims to acquire a clone model $C$ imitating the functionality of victim model $V$. Specifically, adversary seeks to maximize classification accuracy $Acc(C(x; \theta_C), y)$ on the victim model's test set $D_V$, formalized as:

$$\underset{\theta_C}{\operatorname{argmax}} \, E_{(x,y) \sim D_V}[Acc(C(x; \theta_C), y)], \qquad (1)$$

where $\theta_C$ represents the parameter of clone model $C$, and $y$ denotes the ground truth label of test sample $x$.

In practical scenarios, adversaries have limited knowledge of the victim model's architecture, parameters, hyperparameters, and training/testing sets, and can only interact through query APIs. This work adopts a soft-label setting, where the victim model outputs softmax probabilities for user-friendliness. In more realistic settings, adversaries must rely on generative models to synthesize query samples, as they cannot access real samples resembling the defender's private training data. Thus, the clone model $C$'s objective is to minimize the prediction difference between itself and the victim model $V$ on the synthetic query dataset $D_Q$:

$$\underset{\theta_C}{\operatorname{argmin}} \, E_{x \sim D_Q}[d(C(x; \theta_C), V(x; \theta_V))], \qquad (2)$$

where $d(\cdot, \cdot)$ is an indicator measuring difference distance between two objects, and $\theta_V$ is parameters of victim model.

**Defender's Goals and Capabilities.** Defenders lack precise knowledge of the attacker's attack method and duration, but they can use various methods such as anomaly detection and prediction perturbation to resist attacks. The defense objective is twofold: firstly, to impede the adversary's attack, resulting in a clone model with diminished accuracy, and secondly, to
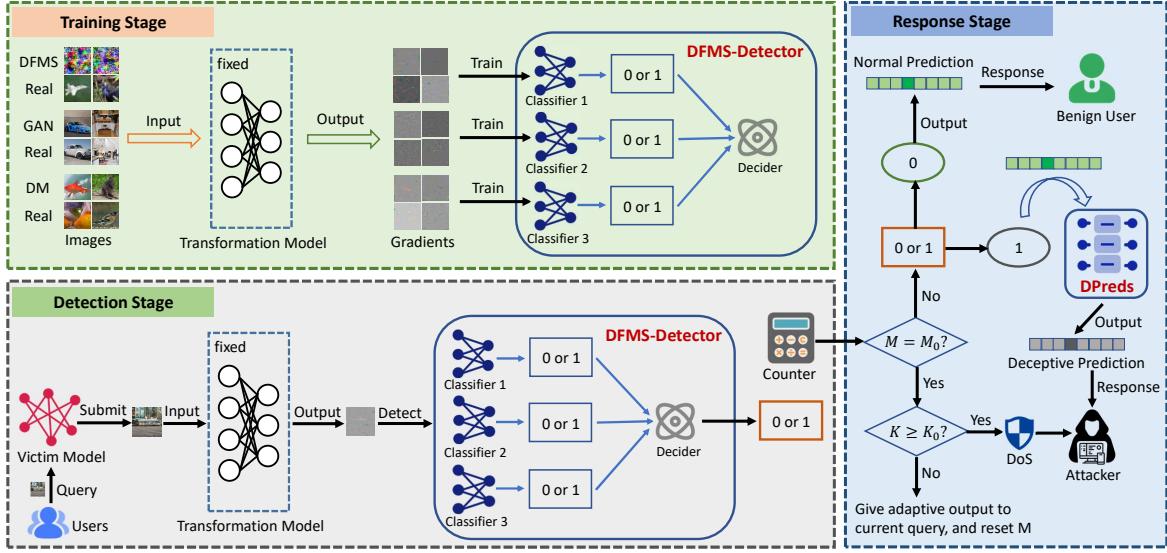
Fig. 1: **Overview of our proposed Model-Guardian.** During training, synthetic images from a randomly selected Data-Free Model Stealing (DFMS) attack, Generative Adversarial Network (GAN), and Diffusion Model (DM) are combined with real images to form three distinct training sets. A pre-trained transformation model converts these data into gradients, which are used to train three sub-detectors, later integrated into a unified detector. In the detection phase, query samples of user are converted into gradients and passed to the detector for evaluation. During the response phase, the model adapts its output based on the query record.

minimize interference with the normal use of benign users. This goal can be formalized as follows:

$$\underset{\theta_V}{\arg\min} E_{(x,y)\sim D_V}[Acc(C(x;\theta_C),y)],$$
$$\text{s.t. } E_{(x,y)\sim D_V}[Acc(V(x;\theta_V),y)] \geq T, \quad (3)$$

where $T$ represents the acceptable threshold for the accuracy of the victim model after degradation.

### B. Overview of Model-Guardian

An overview of our method is shown in Figure 1. Generated images from three sources, along with real images, are converted into gradients. Each sub-detector is trained separately, and their outputs are integrated into an ensemble detector. For each query, the sub-detectors evaluate its maliciousness; if any sub-detector flags the sample as malicious, DFMS-Detector labels it as such. Based on the detection results, benign samples receive normal predictions, while malicious ones trigger the DPreds module for deceptive predictions. Our framework tracks the total queries and malicious samples, and if the proportion of malicious samples in $M_0$ queries exceeds $K_0$, the user's subsequent queries are terminated.

### C. Training DFMS-Detector

**Collecting Training Data.** To improve defense generalization and address common sample synthesis methods, we collect a limited training dataset from three generative sources, each representing a fake class. These sources include data from general data-free model stealing, GANs, and diffusion models. For each synthetic data type, an equal amount of real data is selected as the real class. We focus on the data synthesized by GAN and diffusion model due to attackers' ability, with advancing image generation technology, to create

realistic fake images that closely resemble the victim model's test distribution, potentially bypassing conventional defenses.

**Transforming Images to Gradients.** Our theoretical analysis in Appendix-B and experimental results in Table VII of Appendix-E show that directly training the detector with image-form training data has limited effectiveness and generalization. In line with insights from prior works [24]–[27], we use gradients, a more generalized representation, as the final training data. Gradients filter out image content, retaining only key pixels. Figure 2 illustrates the transformation from image data to gradients, with Class Activation Maps (CAM) extracted from our detector. A pre-trained CNN model is used to convert image data into gradient data, forming a novel training set.

Given a training dataset $D_I = \{(x_i, y_i)\}_{i=1}^{n}$, where $y_i$ represents the label of $x_i$ with two classes: real ($y = 0$) and fake ($y = 1$). We introduce the transformation model $M_T(\cdot)$. The initial step involves inputting $x_i$ into $M_T$ to obtain an output feature vector $u$:

$$u = M_T(x_i). \quad (4)$$

Subsequently, we compute the gradient of $sum(u)$ with respect to input $x_i$:

$$g = \frac{\partial sum(u)}{\partial x_i}, \quad (5)$$

where $sum(\cdot)$ is a summation operation, and $g$ serves as the generalized sample representation.

**Training the Detector.** We normalize the gradients to the range 0–255 for training each sub-detector. Each sub-detector is a binary classifier optimized to distinguish whether an input gradient is from a fake image. The training objective for each
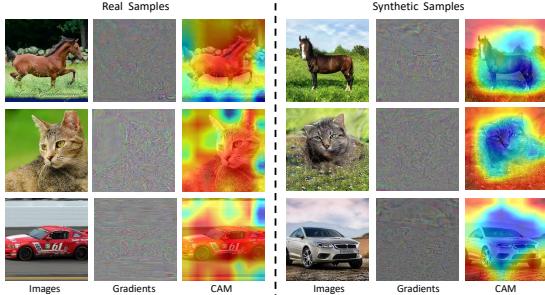
Fig. 2: Visualization of gradients and Class Activate Map (CAM) extracted from detector on real and synthetic images.

sub-detector $F$ is as follows:

$$L = E_{(g,y) \sim D_G}[CE(F(g), y)], \quad (6)$$

where $CE(\cdot, \cdot)$ represents cross-entropy loss, and $D_G$ is the gradient dataset with labels matching the original images. After training, the sub-detectors are combined in parallel to form the ensemble DFMS-Detector.

### D. Detecting Query Sample

**Detecting with Each Sub-Detector.** For user queries, we transform them into gradients and pass them to DFMS-Detector. Each sub-detector independently assesses whether the sample is malicious.

**Ensembling and Making the Final Decision.** A decision-maker evaluates the sample based on outputs of sub-detector. If at least one sub-detector flags the sample as fake ($y = 1$), it is classified as malicious; otherwise, it is benign. A query counter tracks the total number of queries and those identified as malicious.

### E. Returning Adaptive Output

For each query, output strategies depend on the detection results and the number of malicious queries. If the total queries exceed $M_0$ and the proportion of malicious queries exceeds $K_0$, further queries are terminated. Otherwise, adaptive responses are used.

**Normal Predictions for Benign Samples.** If the decision-maker classifies a sample as normal, the original prediction (softmax probability) is returned without modification, ensuring benign users are unaffected.

**Deceptive Predictions for Malicious Samples.** If a query is identified as malicious, the DPreds module is activated to provide deceptive predictions. DPreds is an effective active defense method proposed by us. Adversaries need unaltered probabilities to train clone models, while benign users only require relative class probabilities. Our perturbation algorithm influences malicious queries' outputs while evading detection by adversaries and preserving the order of class labels before and after perturbation.

In the implementation, an identical perturbation value $r$ is added to each class probability of the original prediction $P = $

$[p_1, p_2, ..., p_k]$, followed by normalization to ensure a sum of one as follows:

$$P' = [p_1 + r, p_2 + r, ..., p_k + r],$$
$$P'' = [p_1'', p_2'', ..., p_k''] = \frac{P'}{\sum_{i=1}^{k}(p_i + r)}. \quad (7)$$

The modified prediction result $P''$ is returned to the malicious user, minimizing useful information and hindering effective clone model training.

## IV. EXPERIMENTS

### A. Experimental Setup

**Datasets and Architectures.** In the defense testing phase, we perform experiments on: (1) defense against data-free stealing attacks on CIFAR-10, CIFAR-100, and ImageNet; (2) detection on datasets generated by seven GAN models, including ProGAN [28], StyleGAN [29], StyleGAN2 [30], BigGAN [31], CycleGAN [32], StarGAN [33], and Gau-GAN [34]; (3) detection on datasets generated by six diffusion models, including ADM [35], SD-v1 [36], Guided [35], DALL-E [37], LDM [36], and Glide [38].

We use ResNet-34 as the victim model and ResNet-18 as the clone model. The transformation model is the StyleGAN discriminator pre-trained on LSUN-bedroom, and the classifier for all sub-detectors is a pre-trained ResNet-50 on ImageNet. Details of DFMS-Detector's training set are in Appendix-C.

**Attack Methods.** We evaluate our defense against seven attack methods: DaST [17], MAZE [5], DFME [1], DFMS [4], DisGUIDE [2], Dual Student [6], and QUDA [7]. Additionally, we generate malicious query samples using seven GAN models and six diffusion models to assess detection performance.

**Defense Methods.** We compare our method with eleven defenses: eight active defenses (RS [18], MAD [19], AM [9], EDM [20], NT [21], PoW [22], APGP [10], InI [11]) and three passive defenses (PRADA [12], SEAT [23], FDINET [13]).

**Evaluation Metrics.** (1) Benign Accuracy (BAcc.): percentage of correctly classified benign test data by protected model; (2) Clone Accuracy (CAcc.): percentage of correctly classified benign test data on clone model; (3) Detection Accuracy (DAcc.): percentage of successfully detected malicious samples; (4) False Positive Rate (FPR): percentage of benign samples incorrectly classified as malicious.

**Implementation Details.** When a user's cumulative query count reaches $M_0 = 50k$, Model-Guardian terminates service if malicious queries exceed $K_0 = 50\%$. Adversaries are given query budgets of 20M (CIFAR-10, ImageNet) and 10M (CIFAR-100). Each experiment is run five times, with the average result reported. More details are in the Appendix-C.

### B. Defending against Mainstream Stealing Attacks

**Comparison with Popular Active Defenses.** Table I compares our method with other active defenses against seven attacks on CIFAR-10, CIFAR-100, and ImageNet. Our approach provides strong defense, with clone model accuracies ranging from 10%~15%, 2%~10% and 2%~15% on these datasets, respectively. Importantly, Model-Guardian does not affect the benign accuracy of normal user queries, thanks to the low false

Table I: Performance comparison (%) with popular active defenses.

| Dataset | Method | BAcc. | CAcc. on Testing Attack Method | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | DaST | MAZE | DFME | DFMS | DisGUIDE | Dual Student | QUDA |
| CIFAR-10 | No Defense | 95.54 | 35.18 | 45.60 | 88.10 | 91.24 | 94.02 | 91.36 | 86.94 |
| | RS | 91.81 (-3.73) | 26.39 | 36.42 | 67.84 | 73.90 | 69.08 | 74.41 | 62.69 |
| | MAD | 95.16 (-0.38) | 27.44 | 35.57 | 70.53 | 75.73 | 69.57 | 79.48 | 72.16 |
| | AM | 94.62 (-0.92) | 32.27 | 42.96 | 82.02 | 85.79 | 88.19 | 84.05 | 82.33 |
| | EDM | 95.06 (-0.48) | 32.01 | 42.76 | 83.70 | 83.94 | 89.32 | 87.71 | 83.47 |
| | NT | 94.72 (-0.82) | 23.22 | 30.58 | 61.79 | 71.16 | 68.87 | 63.95 | 60.86 |
| | PoW | 93.22 (-2.32) | 32.96 | 41.83 | 81.76 | 84.22 | 87.63 | 84.96 | 80.68 |
| | APGP | **95.54 (-0.00)** | 21.36 | 26.45 | 52.27 | 58.93 | 61.13 | 59.38 | 55.64 |
| | InI | 94.97 (-0.57) | 31.12 | 41.45 | 78.67 | 83.12 | 86.69 | 83.05 | 76.43 |
| | Model-Guardian | **95.54 (-0.00)** | **10.63** | **11.27** | **11.81** | **12.95** | **14.93** | **12.68** | **13.72** |
| CIFAR-100 | No Defense | 78.52 | 9.48 | 17.81 | 26.46 | 48.83 | 69.47 | 46.52 | 29.35 |
| | RS | 75.53 (-2.99) | 7.58 | 13.39 | 19.53 | 38.19 | 53.74 | 33.02 | 22.89 |
| | MAD | 76.95 (-1.57) | 7.51 | 14.25 | 21.43 | 37.60 | 58.72 | 34.87 | 23.95 |
| | AM | 76.49 (-2.03) | 8.91 | 16.14 | 24.08 | 46.01 | 63.94 | 42.57 | 27.24 |
| | EDM | 77.19 (-1.33) | 8.84 | 16.39 | 25.14 | 44.19 | 63.15 | 42.71 | 26.68 |
| | NT | 77.56 (-0.96) | 7.23 | 11.97 | 18.76 | 33.84 | 50.85 | 32.33 | 21.48 |
| | PoW | 76.71 (-1.81) | 8.66 | 16.13 | 24.24 | 45.80 | 65.72 | 41.89 | 27.39 |
| | APGP | **78.52 (-0.00)** | 5.86 | 11.31 | 15.22 | 29.79 | 42.38 | 29.30 | 18.93 |
| | InI | 78.05 (-0.47) | 8.22 | 16.46 | 23.07 | 44.34 | 64.54 | 40.71 | 26.83 |
| | Model-Guardian | **78.52 (-0.00)** | **2.59** | **3.08** | **5.19** | **8.18** | **9.64** | **7.03** | **7.26** |
| ImageNet | No Defense | 62.08 | 8.86 | 15.47 | 23.84 | 43.69 | 61.64 | 45.88 | 24.97 |
| | RS | 58.87 (-3.21) | 7.06 | 11.72 | 18.69 | 31.78 | 48.29 | 32.69 | 18.58 |
| | MAD | 60.66 (-1.42) | 7.18 | 11.81 | 19.32 | 35.01 | 47.45 | 37.85 | 19.56 |
| | AM | 59.74 (-2.34) | 8.25 | 14.16 | 22.02 | 40.04 | 57.50 | 43.22 | 23.74 |
| | EDM | 60.22 (-1.86) | 8.12 | 14.24 | 22.33 | 40.45 | 57.46 | 42.97 | 23.48 |
| | NT | 61.05 (-1.03) | 6.23 | 10.79 | 16.44 | 31.62 | 43.98 | 34.26 | 18.35 |
| | PoW | 59.44 (-2.64) | 7.96 | 14.68 | 21.86 | 41.37 | 57.68 | 43.14 | 22.94 |
| | APGP | **62.08 (-0.00)** | 5.17 | 9.22 | 14.90 | 24.63 | 38.43 | 29.51 | 14.32 |
| | InI | 61.39 (-0.69) | 7.77 | 13.87 | 22.08 | 38.69 | 54.03 | 42.09 | 22.85 |
| | Model-Guardian | **62.08 (-0.00)** | **2.25** | **2.51** | **4.52** | **10.33** | **14.41** | **11.02** | **4.73** |

Table II: **Performance comparison (%) with popular passive defenses.** To ensure consistency with existing research, we evaluate only the detection performance of our DFMS-Detector. Existing methods assess whether each batch is a sequence of malicious queries. Thus, the results in this table represent the detection accuracy and false positive rate of each method across all batch sequences, indicating the proportion of batch sequences that are identified as malicious.

| Method | bs | DaST | | MAZE | | DFME | | DFMS | | DisGUIDE | | Dual Student | | QUDA | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | DAcc. | FPR | DAcc. | FPR | DAcc. | FPR | DAcc. | FPR | DAcc. | FPR | DAcc. | FPR | DAcc. | FPR |
| PRADA | 50 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| | 500 | 18.00 | 34.00 | 22.00 | 34.00 | 32.00 | 34.00 | 13.00 | 34.00 | 25.00 | 34.00 | 21.00 | 34.00 | 16.00 | 34.00 |
| SEAT | 50 | 13.60 | 0.80 | 10.20 | 0.90 | 11.80 | 0.80 | 9.40 | 0.70 | 13.80 | 0.80 | 12.70 | 0.80 | 11.30 | 0.90 |
| | 500 | 72.00 | 6.00 | 65.00 | 6.00 | 76.00 | 6.00 | 54.00 | 6.00 | 81.00 | 6.00 | 77.00 | 6.00 | 73.00 | 6.00 |
| FDINET | 50 | 86.40 | 1.80 | 89.30 | 1.70 | 89.90 | 1.80 | 89.80 | 1.60 | 91.60 | 1.70 | 92.20 | 1.80 | 94.10 | 1.90 |
| | 500 | 95.00 | 2.00 | 98.00 | 2.00 | 100.00 | 2.00 | 97.00 | 2.00 | 99.00 | 2.00 | 100.00 | 2.00 | 100.00 | 2.00 |
| Model-Guardian | 50 | **97.30** | 0.60 | **98.10** | 0.40 | **99.60** | 0.20 | **99.20** | 0.10 | **98.70** | 0.40 | **98.50** | 0.30 | **97.90** | 0.50 |
| | 500 | **98.00** | 2.00 | **99.00** | 2.00 | **100.00** | 0.00 | **99.00** | 0.00 | **100.00** | 2.00 | **100.00** | 1.00 | 99.00 | 1.00 |

positive rate of the DFMS-Detector. Even with occasional false positives, the class probability relationships in the perturbed output remain unchanged, ensuring correct classification.

**Comparison with Popular Passive Defenses.** We also compare detection performance on CIFAR-10 with three passive defenses. Since the other methods rely on the feature distribution of continuous queries, we vary batch sizes (50 and 500) to assess malicious query behavior. As shown in Table II, our method achieves over 97% detection accuracy and under 2% false positive rate, owing to the strong generalization of our detector, which independently evaluates each sample.

### C. Performance of DFMS-Detector on Malicious Queries from Various Generative Models

In Table III and Table IV, our DFMS-Detector demonstrates superior detection performance on all images generated by various GANs and diffusion models, with average detection accuracies of 98.43% and 97.00%, respectively.

### D. Effectiveness of Deceptive Predictions

To evaluate DPreds' impact on model stealing defense, we test it independently against six data-free stealing attacks. As shown in Table V with a perturbation value of 0.5, the module reduces clone model accuracy by 19.24%~34.57% across various attack methods. For a clearer view of the perturbation algorithm's effect, we visualize the class probability

Table III: **Detection performance (%) of DFMS-Detector on various GAN-generated images.** The meanings of each metric are consistent with Table II, and the batch size is 50.

| Method | ProGAN | | StyleGAN | | StyleGAN2 | | BigGAN | | CycleGAN | | StarGAN | | GauGAN | | Total Avg. | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | DAcc. | FPR | DAcc. | FPR | DAcc. | FPR | DAcc. | FPR | DAcc. | FPR | DAcc. | FPR | DAcc. | FPR | DAcc. | FPR |
| PRADA | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| SEAT | 9.00 | 3.00 | 12.00 | 4.00 | 10.00 | 4.00 | 18.00 | 3.00 | 13.00 | 5.00 | 21.00 | 2.00 | 24.00 | 5.00 | 15.29 | 3.71 |
| FDINet | 91.00 | 0.00 | 89.00 | 0.00 | 86.00 | 0.00 | 93.00 | 2.00 | 89.00 | 4.00 | 95.00 | 0.00 | 92.00 | 1.00 | 90.71 | 1.00 |
| Model-Guardian | **100.00** | 0.00 | **97.00** | 1.00 | **96.00** | 2.00 | **99.00** | 0.00 | **99.00** | 2.00 | **100.00** | 0.00 | **98.00** | 1.00 | **98.43** | 0.86 |

Table IV: **Detection performance (%) of DFMS-Detector on various diffusion-generated images.** The meanings of each metric are consistent with Table II, and the batch size is 50.

| Method | ADM | | SD-v1 | | Guided | | DALL-E | | LDM | | Glide | | Total Avg. | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | DAcc. | FPR | DAcc. | FPR | DAcc. | FPR | DAcc. | FPR | DAcc. | FPR | DAcc. | FPR | DAcc. | FPR |
| PRADA | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| SEAT | 5.00 | 3.00 | 6.00 | 3.00 | 8.00 | 2.00 | 3.00 | 0.00 | 3.00 | 1.00 | 5.00 | 0.00 | 5.00 | 1.50 |
| FDINet | 82.00 | 2.00 | 80.00 | 1.00 | 86.00 | 1.00 | 79.00 | 1.00 | 75.00 | 1.00 | 84.00 | 0.00 | 81.00 | 1.00 |
| Model-Guardian | **100.00** | 0.00 | **94.00** | 2.00 | **97.00** | 1.00 | **95.00** | 1.00 | **98.00** | 1.00 | **98.00** | 1.00 | **97.00** | 1.00 |

Table V: **Effectiveness (%) of Deceptive Predictions.** The results indicate the accuracy of clone model obtained by attacker in two situations: no defense and with only DPreds.

| Method | DaST | MAZE | DFME | DFMS | DisGUIDE | QUDA |
|---|---|---|---|---|---|---|
| No Defense | 35.18 | 45.60 | 88.10 | 91.24 | 94.02 | 86.94 |
| w/ DPreds | **15.88** | **26.36** | **53.53** | **62.07** | **61.15** | **53.03** |
| Decline rate | 19.30 | 19.24 | 34.57 | 29.17 | 32.87 | 33.91 |



Fig. 3: Visualization of class probabilities before perturbation (the second and fifth columns) and after perturbation (the third and last columns) on synthetic query images from six different sources (DaST, StyleGAN, ADM located in the first column and DFME, GauGAN, DALL-E located in the third column).

distribution before and after perturbation in Figure 3. The visualization shows that while the class relationships remain intact, the probability values change, creating smoother peaks that mislead adversaries and lower clone model accuracy. Normal users remain unaffected, as they rely only on top-1 classification or the order of probabilities. Additionally, we provide t-SNE visualizations of the prediction distribution on CIFAR-10 before and after perturbations in the Appendix-D.

### E. Other Experiments and Discussions

We also explore the impact of the transformation model, training data size, and perturbation value on our method, with results in Appendix-E. Discussions on hyperparameter selection, computational load, multi-adversary collusion, and extensibility are in Appendix-F.

## V. CONCLUSION

In this paper, we propose Model-Guardian, a defense against data-free model stealing attacks. By collecting synthetic and real data from three sources, we transform images into gradients, a more generalized representation. We then train binary classifiers and integrate them into the DFMS-Detector for malicious query detection. Additionally, we introduce the Deceptive Predictions algorithm to perturb outputs for malicious users. Extensive experiments show Model-Guardian outperforms existing defenses, effectively detecting covert malicious queries. Our work aims to inspire further research in defending against data-free model stealing attacks.

## REFERENCES

[1] Jean-Baptiste Truong et al., "Data-free model extraction," in *CVPR*, 2021, pp. 4771–4780.
[2] Jonathan Rosenthal et al., "Disguide: Disagreement-guided data-free model extraction," *AAAI*, 2023.
[3] Yunfei Yang et al., "Stms: An out-of-distribution model stealing method based on causality," in *IJCNN*, 2024, pp. 1–8.
[4] Sunandini Sanyal et al., "Towards data-free model stealing in a hard label setting," in *CVPR*, 2022, pp. 15284–15293.
[5] Sanjay Kariyappa et al., "Maze: Data-free model stealing attack using zeroth-order gradient estimation," in *CVPR*, 2021, pp. 13814–13823.
[6] James Beetham et al., "Dual student networks for data-free model stealing," in *ICLR*, 2022.
[7] Zijun Lin et al., "Quda: Query-limited data-free model extraction," in *AsiaCCS*, 2023, pp. 913–924.
[8] Yunfei Yang et al., "Dualcos: Query-efficient data-free model stealing with dual clone networks and optimal samples," in *ICME*, 2024.
[9] Sanjay Kariyappa et al., "Defending against model stealing attacks with adaptive misinformation," in *CVPR*, 2020, pp. 770–778.
[10] Anda Cheng et al., "Apgp: Accuracy-preserving generative perturbation for defending against model cloning attacks," in *ICASSP*, 2023, pp. 1–5.
[11] Jun Guo et al., "Isolation and induction: Training robust deep neural networks against model stealing attacks," in *ACM MM*, 2023.
[12] Mika Juuti et al., "Prada: protecting against dnn model stealing attacks," in *EuroS&P*. IEEE, 2019, pp. 512–527.
[13] Hongwei Yao et al., "Fdinet: Protecting against dnn model extraction via feature distortion index," *arXiv preprint arXiv:2306.11338*, 2023.
[14] Nicolas Papernot et al., "Practical black-box attacks against machine learning," in *AsiaCCS*, 2017, pp. 506–519.

[15] Tribhuvanesh Orekondy et al., "Knockoff nets: Stealing functionality of black-box models," in *CVPR*, 2019, pp. 4954–4963.

[16] Yaxin Xiao et al., "Mexmi: Pool-based active model extraction crossover membership inference," *NIPS*, vol. 35, pp. 10203–10216, 2022.

[17] Mingyi Zhou et al., "Dast: Data-free substitute training for adversarial attacks," in *CVPR*, 2020, pp. 234–243.

[18] Taesung Lee et al., "Defending against neural network model stealing attacks using deceptive perturbations," in *S&P Workshop*. IEEE, 2019.

[19] Tribhuvanesh Orekondy et al., "Prediction poisoning: Towards defenses against dnn model stealing attacks," in *ICLR*, 2020.

[20] Sanjay Kariyappa et al., "Protecting dnns from theft using an ensemble of diverse models," in *ICLR*, 2021.

[21] Haoyu Ma et al., "Undistillable: Making a nasty teacher that cannot teach students," in *ICLR*, 2021.

[22] Adam Dziedzic et al., "Increasing the cost of model extraction with calibrated proof of work," in *ICLR*, 2022.

[23] Zhanyuan Zhang et al., "Seat: similarity encoder by adversarial training for detecting model extraction attack queries," in *CCS Workshop*, 2021.

[24] Sheng-Yu Wang et al., "Cnn-generated images are surprisingly easy to spot... for now," in *CVPR*, 2020, pp. 8695–8704.

[25] Yiming Li et al., "Defending against model stealing via verifying embedded external features," in *AAAI*, 2022, vol. 36, pp. 1464–1472.

[26] Yun Liu et al., "Detection of gan generated image using color gradient representation," *JVCIR*, vol. 95, pp. 103876, 2023.

[27] Chuangchuang Tan et al., "Learning on gradients: Generalized artifacts representation for gan-generated images detection," in *CVPR*, 2023, pp. 12105–12114.

[28] Tero Karras et al., "Progressive growing of gans for improved quality, stability, and variation," in *ICLR*, 2018.

[29] Tero Karras et al., "A style-based generator architecture for generative adversarial networks," in *CVPR*, 2019, pp. 4401–4410.

[30] Tero Karras et al., "Analyzing and improving the image quality of stylegan," in *CVPR*, 2020, pp. 8110–8119.

[31] Andrew Brock, Jeff Donahue, and Karen Simonyan, "Large scale gan training for high fidelity natural image synthesis," in *ICLR*, 2018.

[32] Jun-Yan Zhu et al., "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *ICCV*, 2017, pp. 2223–2232.

[33] Yunjey Choi et al., "Stargan: Unified generative adversarial networks for multi-domain image-to-image translation," in *CVPR*, 2018.

[34] Taesung Park et al., "Semantic image synthesis with spatially-adaptive normalization," in *CVPR*, 2019, pp. 2337–2346.

[35] Prafulla Dhariwal et al., "Diffusion models beat gans on image synthesis," *NIPS*, vol. 34, pp. 8780–8794, 2021.

[36] Robin Rombach et al., "High-resolution image synthesis with latent diffusion models," in *CVPR*, 2022, pp. 10684–10695.

[37] Aditya Ramesh et al., "Zero-shot text-to-image generation," in *ICML*. PMLR, 2021, pp. 8821–8831.

[38] Alexander Quinn Nichol et al., "Glide: Towards photorealistic image generation and editing with text-guided diffusion models," in *ICML*. PMLR, 2022, pp. 16784–16804.

[39] Xu Zhang et al., "Detecting and simulating artifacts in gan fake images," in *WIFS*. IEEE, 2019, pp. 1–6.

[40] Joel Frank et al., "Leveraging frequency analysis for deep fake image recognition," in *ICML*. PMLR, 2020, pp. 3247–3258.

[41] Ricard Durall et al., "Watch your up-convolution: Cnn based generative deep neural networks are failing to reproduce spectral distributions," in *CVPR*, 2020, pp. 7890–7899.

[42] Chuangchuang Tan et al., "Rethinking the up-sampling operations in cnn-based generative network for generalizable deepfake detection," in *CVPR*, 2024, pp. 28130–28139.

## APPENDIX

### A. The Workflow of Model Stealing Attack

As shown in Figure 4, model stealing attacks involve adversaries continuously querying APIs with crafted samples to acquire annotated data for training a functionally equivalent clone model. Beyond utilizing the stolen model without cost, attackers can employ it for adversarial attacks, membership inference attacks, and model inversion attacks.
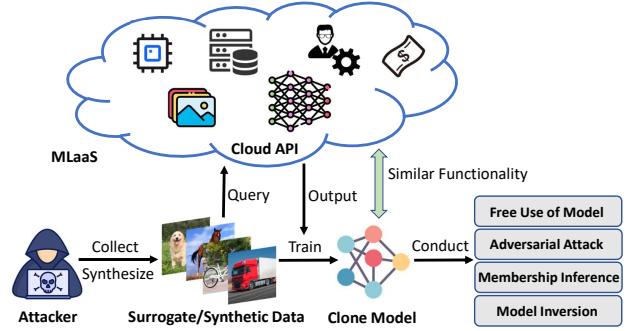


Fig. 4: **Model stealing attack and its vulnerabilities.** Attackers will query models deployed in the cloud through APIs using surrogate or synthetic data to obtain corresponding predictions. They can then use these annotated data to train a clone model with similar functionality to the original models and carry out further malicious actions.

### B. Insights on Defeating Data-Free Stealing

**Fundamental Characteristic of Synthetic Images: Artifacts.** In the context of Machine Learning as a Service (MLaaS), benign users query target models with natural, real data. However, a significant distinction of data-free model stealing attacks from previous attacks is the absence of real data, as they rely on synthetic data generated by generative models. Based on this fundamental difference, a straightforward and feasible defense against such attacks is to identify the synthetic queries used by attackers. Current data synthesis techniques, including generative adversarial networks (GANs) and diffusion models, produce data through generative processes. Both GANs and diffusion models involve numerous up-sampling operations in their workflows. Research in the field of deepfake detection [39]–[42] has demonstrated that due to these up-sampling components, synthetic data generated by these models inherently carry unnatural features, known as artifacts. Artifacts refer to unnatural, artificially processed traces, or regions. Consequently, artifacts can serve as a reliable criterion for distinguishing between normal and malicious queries when defending against data-free stealing.

**Enhancing Defense Generalization: Transforming Data Dependency to Model Dependency.** Building on the previous analysis, a straightforward method for detecting malicious samples involves collecting real and synthetic images to train a binary classifier as a detector. However, prior research [24] has demonstrated that such detectors suffer from poor generalization, effectively identifying images generated by the same generator used during training but failing with images from other generators. This issue arises due to the detector's over-reliance on training data. To address this challenge, inspired by previous work [24]–[27], we propose using a more generalized representation, namely, the data gradients on a pre-trained model, as the detector's training data. Since these gradients retain only the essential pixels relevant to the pre-trained model's target task, they exclude the content of images.

Table VI: Composition of the training dataset for DFMS-Detector.

| Dataset Name | Real Source | Fake Source | # of images (real/fake) |
|---|---|---|---|
| Data Subset 1 | TinyImageNet | DFME | 20k/20k |
| Data Subset 2 | LSUN | ProGAN | 20k/20k |
| Data Subset 3 | ImageNet | ADM | 20k/20k |

The remaining core pixels help the detector learn the common characteristics of this data type. Consequently, we transform the original data dependency problem into a model dependency problem. By converting synthetic data from different sources into gradient data using the same pre-trained model, we can reliably detect artifacts, ensuring the detector's effectiveness across various types of synthetic data.

**Defense vs. Utility: A Comprehensive Defense Approach.** The aforementioned detector serves as just one component of our defense against data-free stealing attacks. Our ultimate objective is to prevent attackers from obtaining a well-performing clone model while ensuring that legitimate users can continue to use the protected model without disruption. To achieve this goal, we must consider multiple factors, including: (1) enhancing the reliability of detection results to minimize false positives, (2) managing malicious query samples and malicious users, and (3) maximizing user experience for benign users. Addressing these challenges necessitates the design of a comprehensive defense framework.

### C. Other Experimental Setup

**Composition of the Training Dataset for DFMS-Detector.** In constructing the training set for DFMS-Detector, we collect a representative subset of samples from each of three different sources, which is a viable strategy for defenders. In detail, We simulate the data-free stealing process of DFME to query the defensive model, obtain a trained generator, and then input random noise to it to synthesize the synthetic samples of our first type of training data. The second type of synthetic training data is generated by a currently mainstream GAN model called ProGAN. The third type of synthetic data is synthesized using the popular diffusion model ADM. For each type of synthetic data, as samples with the class fake in the corresponding training subset, we additionally collect an equal amount of real samples with semantic class relevance to the synthetic data as samples with the class real in this training subset. The complete information regarding the training set is provided in Table VI.

**More Implementation Details.** During DFMS-Detector training, the Adam optimizer is employed with an initial learning rate of $5 \times 10^{-4}$. Decay involves multiplying the learning rate by 0.9 every ten epochs, with a total of 200 training epochs and a batch size of 32. Data augmentation methods, including random cropping, horizontal flipping, blurring, and compression, are utilized. The trained model serves as a classifier for the detector.
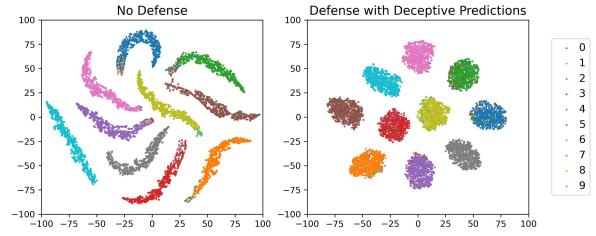


Fig. 5: Visualization of t-SNEs for both normal and defensive ResNet-34 on CIFAR-10. Each dot represents one data point.

Table VII: **Ablation study (%) on the contributions of different modules.** The results indicate the accuracy of clone model obtained by attacker using various data-free model stealing attacks under different defense settings.

| Method | MAZE | DFME | DFMS | DisGUIDE | QUDA |
|---|---|---|---|---|---|
| No Defense | 45.60 | 88.10 | 91.24 | 94.02 | 86.94 |
| Model-Guardian | **11.27** | **11.81** | **12.95** | **14.93** | **13.72** |
| w/o Gradient Trans. | 23.74 | 13.58 | 48.63 | 45.32 | 37.05 |
| w/o DFMS-Detector | 26.36 | 53.53 | 62.07 | 61.15 | 53.03 |
| w/o DPreds | 13.95 | 15.22 | 16.86 | 17.94 | 18.87 |

**Experimental Platform.** All our experiments are conducted on a server running the 64-bit Ubuntu 20.04.6 operating system. The server is equipped with an Intel(R) Xeon(R) Silver 4314 CPU @ 2.40GHz, 504GB of memory, and one NVIDIA A100 PCIe GPU with 40GB of memory. In the specific implementation of the source code, we use Python 3.10.4, Pytorch 1.11.0 and CUDA 11.8.

### D. T-SNE Visualization for the Prediction Distribution of Victim Model on CIFAR-10

We illustrate t-SNE visualization for the prediction distribution of the victim model on CIFAR-10 before and after perturbations in Figure 5. It is evident that after perturbations, although the model maintains normal classification, its output distribution undergoes a significant shift, deceiving adversaries' attempts to clone the victim model by mimicking its predictions.

### E. More Ablation Studies

In this subsection, we conduct some ablation studies to further understand the functions of various modules of Model-Guardian, the impact of transformation models on the performance, the impact of the number of detector's training data on the performance, and the effectiveness of perturbation algorithm under different perturbation values.

**Contributions of Different Modules.** Table VII demonstrates that training the detector directly with images, without gradient transformation, results in reduced defense efficacy compared to the complete Model-Guardian. Particularly noteworthy is the significant improvement in effectiveness when DPreds module is employed for defense alone, compared to a scenario with no defense, underscoring the efficacy of the perturbation algorithm. Additionally, excluding the perturbation algorithm while retaining only the DFMS-Detector leads to

Table VIII: **Detection performance (%) of DFMS-Detector with different transformation models.** Since only the detection performance of our DFMS-Detector is evaluated here, the results in this table directly represent the overall detection accuracy and false positive rate on all test samples.

| Transformation Model | DaST | | MAZE | | DFME | | DFMS | | DisGUIDE | | Dual Student | | QUDA | | Total Avg. | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | DAcc. | FPR | DAcc. | FPR | DAcc. | FPR | DAcc. | FPR | DAcc. | FPR | DAcc. | FPR | DAcc. | FPR | DAcc. | FPR |
| Input Image | 64.82 | 8.34 | 79.18 | 7.12 | 97.43 | 4.59 | 61.09 | 6.56 | 76.34 | 9.59 | 71.95 | 9.78 | 56.92 | 9.88 | 72.53 | 7.98 |
| ResNet18 | 73.06 | 7.28 | 85.22 | 4.66 | 96.38 | 3.06 | 69.28 | 4.63 | 83.12 | 6.84 | 80.02 | 7.79 | 65.13 | 6.43 | 78.89 | 5.81 |
| CLIP-ResNet50 | 81.12 | 6.24 | 93.72 | 5.03 | 98.07 | 3.24 | 76.04 | 4.18 | 91.59 | 6.07 | 87.43 | 5.13 | 74.95 | 4.32 | 86.13 | 4.89 |
| ProGAN-Discriminator | 83.35 | 7.16 | 95.46 | 6.05 | 98.12 | 2.90 | 85.77 | 3.23 | 94.03 | 5.34 | 92.75 | 5.10 | 78.62 | 4.19 | 89.73 | 4.85 |
| StyleGAN-Discriminator | **91.47** | 6.05 | **97.16** | 5.23 | **98.36** | 2.32 | **93.75** | 3.09 | **96.37** | 5.01 | **95.58** | 3.65 | **91.74** | 3.52 | **94.92** | 4.12 |

Table IX: **Detection performance (%) of DFMS-Detector under different numbers of training data.** The meanings of each metric are consistent with Table VIII.

| Num. of Train. Data | DaST | | MAZE | | DFME | | DFMS | | DisGUIDE | | Dual Student | | QUDA | | Total Avg. | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | DAcc. | FPR | DAcc. | FPR | DAcc. | FPR | DAcc. | FPR | DAcc. | FPR | DAcc. | FPR | DAcc. | FPR | DAcc. | FPR |
| 5k | 85.97 | 8.13 | 91.17 | 9.92 | 90.92 | 8.35 | 90.38 | 7.56 | 92.33 | 9.97 | 89.49 | 6.93 | 85.22 | 8.49 | 89.35 | 8.48 |
| 10k | 87.28 | 7.40 | 94.32 | 8.68 | 94.66 | 5.77 | 92.02 | 4.87 | 93.05 | 9.43 | 92.83 | 7.45 | 89.35 | 6.06 | 91.93 | 7.09 |
| 20k | 88.92 | 6.56 | 96.05 | 6.02 | 95.97 | 3.38 | 91.80 | 5.64 | 95.14 | 6.72 | 95.26 | 5.24 | 91.08 | 5.43 | 93.46 | 5.57 |
| 40k | 91.47 | 6.05 | 97.16 | 5.23 | **98.36** | 2.32 | **93.75** | 3.09 | **96.37** | 5.01 | 95.58 | 3.65 | **91.74** | 3.52 | **94.92** | 4.12 |
| 60k | **91.83** | 6.22 | **97.21** | 5.39 | 98.19 | 2.54 | 93.14 | 3.71 | 95.84 | 5.32 | **96.02** | 3.51 | 91.59 | 4.12 | 94.83 | 4.40 |

Table X: **Active defense effects (%) of DPreds under different perturbation values.** The results indicate the accuracy of clone model obtained by attacker using various data-free model stealing attacks under defense settings of different perturbation values.

| Perturbation Value | DaST | MAZE | DFME | DFMS | DisGUIDE | Dual Student | QUDA |
|---|---|---|---|---|---|---|---|
| 0.00 | 35.18 | 45.60 | 88.10 | 91.24 | 94.02 | 91.36 | 86.94 |
| 0.05 | 19.04 | 24.41 | 71.15 | 75.73 | 79.93 | 74.82 | 70.44 |
| 0.10 | 16.26 | 22.90 | 72.34 | 74.52 | 77.04 | 71.90 | 68.56 |
| 0.20 | **14.71** | **21.23** | 66.04 | 70.59 | 72.52 | 67.35 | 65.21 |
| 0.50 | 15.88 | 26.36 | **53.53** | 62.07 | 61.15 | 57.56 | **53.03** |
| 1.00 | 15.24 | 26.56 | 56.15 | 59.25 | **60.58** | 56.83 | 53.26 |
| 2.00 | 14.95 | 23.35 | 53.60 | **58.42** | 61.87 | **56.04** | 54.61 |

a minor decrease in defense performance. This highlights the high detection accuracy of our detector, capable of promptly terminating adversary access to the service.

**Detection Performance of DFMS-Detector with Different Transformation Models.** The transformation model plays an important role in implementing gradient transformation in our method, and here we explore the impact of its selection on the performance of the DFMS-Detector. Specifically, we select different types of models such as classification model, contrastive learning model, and discriminator of GAN as transformation models. We can draw three intuitive conclusions from the comparison results in Table VIII: (a) converting the training data of the detector from the input image to gradient data can improve the generalization ability of the detector, because gradient transformation filters out complex image content while retaining distinguishable key pixels; (b) The performance of GAN's discriminator is better than that of classification model and contrastive learning model, which reflects the natural advantage of GAN's discriminator in identifying artifacts in synthesized images caused by the up-sampling operation of the generative model; (c) The effectiveness of discriminators with different structures may also vary. These conclusions indicate that gradient transformation can

transform the original data-dependency problem into a model-dependency problem, greatly improving the generalization of the DFMS-Detector. Given the best performance of the discriminator of StyleGAN, we will apply it as the final transformation model to our framework.

**Detection Performance of DFMS-Detector under Different Numbers of Training Data.** In order to test the impact of different numbers of training data on the performance of the proposed DFMS-Detector, we collect training data of 5k, 10k, 20k, 40k, and 60k, with each setting having half the data number of the synthetic and real samples. During testing, we use 10k synthetic data and 10k real data each. Table IX demonstrates that as the number of training data increases, the detection performance of DFMS-Detector is better, but when it reaches above 40k, the performance increases very little and overfitting occurs. Therefore, the training data number of the DFMS-Detector we used in other experiments is fixed at 40k.

**Active Defense Effects of DPreds under Different Perturbation Values.** We also explore the impact of different perturbation values on the effectiveness of our active defense method DPreds. In order to facilitate direct observation and obtain experimental results, we remove the DFMS-Detector and only retain the active defense module DPreds. From Table X, it

can be seen that as the perturbation value increases, the defense effect will gradually strengthen. In most attack methods, the defense effect reaches its peak when the perturbation value reaches 0.20 or 0.50. Even if the perturbation value continues to increase, the corresponding defense effect does not show a significant increasing trend, but gradually stabilizes.

### F. Discussion

**Selection of Thresholds $M_0$ and $K_0$.** $M_0$ and $K_0$ are critical hyperparameters in our approach, significantly influencing the overall performance of Model-Guardian. A larger $M_0$ means that the user's query count needs to accumulate to a higher number before deciding whether to terminate the service, potentially allowing the attacker to obtain a more accurate clone model. Conversely, a smaller $M_0$ can terminate malicious queries promptly but increases the likelihood of falsely terminating services for benign users. Similarly, a larger $K_0$ means more samples must be deemed malicious before service termination, reducing the detector's sensitivity and allowing attackers to evade detection. A smaller $K_0$ increases the detector's sensitivity but also raises the chance of false positives for benign users. Empirically, we set $M_0$ to $50k$ and $K_0$ to 50% in our experiments, which we believe to be a reasonable combination, balancing the protection against attackers and the normal usage for benign users.

**Computational Cost.** Our method requires each input sample to be transformed into a gradient image and uses a perturbation algorithm to modify the predictions of samples deemed malicious, introducing some computational overhead. We conduct a precise analysis of this computational cost. Specifically, converting a high-resolution $256 \times 256$ image to a gradient image takes only $7.70 \times 10^{-2}$ seconds, and perturbing a 10-class probability vector takes just $1.65 \times 10^{-5}$ seconds. These results indicate that the additional overhead introduced by our defense method is minimal, thus not limiting its practicality.

**Collusion Attacks with Multiple Accounts.** To evade our detection, attackers might resort to registering multiple accounts or collaborating with others to conduct a collusion attack, where multiple accounts or attackers query the victim model to train a single clone model. Our method's inability to resist such collusion attacks is a limitation. However, collusion attacks require higher costs compared to ordinary attacks, and our DPreds module will continue to hinder the training of the clone model to some extent in this scenario, thereby reducing the effectiveness of collusion attacks.

**Flexibility and Extensibility.** In this paper, we focus only on the image classification task. Future research could explore extending our method to defend against data-free model stealing attacks in other types of tasks. Given our framework's flexibility and ease of recombination, it can be adapted with minor modifications to apply to tasks such as object detection, semantic segmentation, image processing, natural language processing, and speech recognition. Additionally, replacing the DFMS-Detector component in Model-Guardian with other suitable anomaly detection algorithms can effectively defend against other types of model stealing attacks based on adversarial data or surrogate data. Overall, our method offers considerable flexibility and extensibility.