

DEEPTASTER: Adversarial Perturbation-Based Fingerprinting to Identify Proprietary Dataset Use in Deep Neural Networks

Seonhye Park¹, Alsharif Abuadbbba², Shuo Wang², Kristen Moore², Yansong Gao²,

Hyoungshick Kim¹, Surya Nepal²,

¹Sungkyunkwan University, Republic of Korea

{qkrtjsgp08,hyoung}@skku.edu

²CSIRO Data61, Australia

{sharif.abuadbbba,shuo.wang,kristen.moore,garrison.gao,surya.nepal}@data61.csiro.au

ABSTRACT

Training deep neural networks (DNNs) requires large datasets and powerful computing resources, which has led some owners to restrict redistribution without permission. Watermarking techniques that embed confidential data into DNNs have been used to protect ownership, but these can degrade model performance and are vulnerable to watermark removal attacks. Recently, DEEPJUDGE was introduced as an alternative approach to measuring the similarity between a suspect and a victim model. While DEEPJUDGE shows promise in addressing the shortcomings of watermarking, it primarily addresses situations where the suspect model copies the victim’s architecture. In this study, we introduce DEEPTASTER, a novel DNN fingerprinting technique, to address scenarios where a victim’s data is unlawfully used to build a suspect model. DEEPTASTER can effectively identify such DNN model theft attacks, even when the suspect model’s architecture deviates from the victim’s. To accomplish this, DEEPTASTER generates adversarial images with perturbations, transforms them into the Fourier frequency domain, and uses these transformed images to identify the dataset used in a suspect model. The underlying premise is that adversarial images can capture the unique characteristics of DNNs built with a specific dataset. To demonstrate the effectiveness of DEEPTASTER, we evaluated the effectiveness of DEEPTASTER by assessing its detection accuracy on three datasets (CIFAR10, MNIST, and Tiny-ImageNet) across three model architectures (ResNet18, VGG16, and DenseNet161). We conducted experiments under various attack scenarios, including transfer learning, pruning, fine-tuning, and data augmentation. Specifically, in the Multi-Architecture Attack scenario, DEEPTASTER was able to identify all the stolen cases across all datasets, while DEEPJUDGE failed to detect any of the cases.

CCS CONCEPTS

• Security and privacy; • Computing methodologies → Machine learning; Artificial intelligence;

KEYWORDS

Neural networks; Model theft; DNN Fingerprinting

1 INTRODUCTION

Organisations looking to commercialise their proprietary Deep Neural Network (DNN) models via Machine Learning as a Service (MLaaS) must be wary of the potential security risks entailed, notably the potential theft of DNN models or datasets [37]. DNN models are constructed on huge datasets, which are meticulously

collected, processed, organised, and labelled, usually requiring significant expense. Therefore, ensuring the protection of model and dataset ownership becomes crucial.

We need to consider several attack scenarios related to the model and data theft. A recent data breach at “Capital One” showed the risk of insider attacks on cloud platforms, where someone with unauthorised access could steal data stored on the cloud server [31]. This security incident highlights the potential misuse of datasets by insider attackers, who could covertly steal a proprietary dataset and incorporate it into their own DNN models without the owner’s consent. Another risk is external attackers stealing a DNN model by querying it through MLaaS APIs. Recent studies (e.g., [33, 40, 46]) have shown that DNN model theft attacks can be effectively conducted, even within real-world services.

Existing DNN intellectual property (IP) protection mechanisms fall into two categories: *DNN watermarking* and *DNN fingerprinting*. DNN watermarking involves embedding the owner’s information (i.e., a watermark) into a proprietary model [1, 2, 6, 7, 9, 20, 38, 41, 47]. Model ownership can be verified by extracting an identical or similar watermark from a suspected model. There have been many proposals for developing effective DNN watermarking schemes. However, DNN watermarking presents two inherent limitations: (a) DNN watermarking is invasive by design, as it necessitates modifications to the original DNN model to embed a watermark, potentially altering the model’s behaviour [42, 47]. (b) DNN watermarking lacks sufficient resilience to adversarial attacks [32, 45]. Aiken *et al.* [3] showed that attackers could effectively manipulate neurons or channels in DNN layers that contribute to the embedded watermark for most state-of-the-art DNN watermarking schemes. Lukas *et al.* [27] recently demonstrated that transfer learning could remove nearly all of the tested 11 watermarking schemes.

In contrast, DNN fingerprinting is *non-invasive* by design as it leverages the unique characteristics (i.e., fingerprinting features) of each DNN model without modifying the model itself. A verifier can identify the model by examining these fingerprinting features [5, 28]. Generally, a single fingerprinting feature is insufficient to identify a model constructed through model theft and adaptive attacks [8]. Chen *et al.* [8] recently introduced DEEPJUDGE, a state-of-the-art fingerprinting scheme that leverages multiple fingerprinting features to protect a model’s copyright. However, as DEEPJUDGE utilises fingerprinting features tied to the model’s parameters, it may struggle to detect unauthorised use of the protected training dataset if a suspect DNN model comprises different parameters or uses a distinct model architecture. Moreover, DEEPJUDGE is not sufficiently effective in detecting models constructed through *transfer*

learning [39]. Our experimental results indicate that DEEPJUDGE’s detection accuracy is significantly degraded for models constructed through transfer learning. Consequently, DEEPJUDGE may fail to identify instances where a victim’s data is illicitly used to construct a suspect model with an architecture that differs from the original model of the victim. To tackle such attack scenarios, we propose a novel DNN fingerprinting scheme dubbed DEEPTASTER.

In this paper, we demonstrate that the spectra of gradient-based adversarial examples can be used to identify the characteristics of a dataset used to train a DNN model, particularly regarding the model’s decision boundaries. The adversarial perturbation images generated by gradient-based attacks preserve both the dataset and model characteristics, indicating a commonality among models trained on the same dataset. Our empirical analysis reveals that adversarial images’ characteristics are more distinctively conserved in the Discrete Fourier Transform (DFT) domain compared to the spatial domain. Adversarial images typically contain more noise than standard images, a consequence of the adversarial perturbation process. We have observed that these noises are more noticeable in the frequency and spatial domains. As depicted in Appendix A, adversarial examples derived from three different architectures (ResNet18, VGG16, and Densenet161) on identical datasets (CIFAR10, MNIST, or Tiny-ImageNet) display significant similarity within the DFT domain. In contrast, adversarial examples within the spatial domain appear completely blacked out, making them visually indistinguishable across datasets. Inspired by these findings, we introduce DEEPTASTER, a method for detecting DNN model theft attacks. DEEPTASTER generates multiple adversarial images with perturbations, transforms them into the DFT domain, and uses their statistical properties as features to identify the dataset used to train a suspect model. Our experimental findings confirm that DEEPTASTER can accurately identify the dataset used to construct a suspect model, given that the architecture of the suspect model is known in advance. Our key contributions are summarised as follows:

- We introduce DEEPTASTER, a novel DNN fingerprinting method designed to identify known model architectures trained on stolen datasets. DEEPTASTER generates adversarial images, transforms them into the DFT domain, and uses these transformed images to discern the unique characteristics of the dataset used to train a suspect model.
- We evaluate the resilience of DEEPTASTER against eight attack scenarios, including multi-architectures, data augmentation, re-training, transfer learning, fine-tuning, pruning, transfer learning with data augmentation, and transfer learning with pretrained model. These evaluations are conducted across three datasets (CIFAR10, MNIST, and Tiny-ImageNet) and three model architectures (ResNet18, VGG16, and DenseNet161). Our experimental results indicate that DEEPTASTER considerably outperforms DEEPJUDGE in most scenarios. For example, in the Multi-Architecture Attack scenario, DEEPTASTER successfully detected all the stolen cases across all datasets. In contrast, DEEPJUDGE failed to detect four attack cases including transfer learning and data augmentation for the CIFAR10 dataset.

2 BACKGROUND

2.1 Adversarial Perturbation and Attack

An adversarial perturbation refers to an intentionally created perturbation of an input sample that can lead to its misclassification by a machine learning model [12, 30, 44]. Gradient-based adversarial attacks are well-known for generating such perturbations, including the fast gradient sign method (FGSM) [12]. Gradient-based adversarial attacks generate a minimal perturbation to the input sample in a direction that most significantly impacts the prediction of the target classifier. This “small modification,” which might be as subtle as changing a single pixel’s color, can potentially disrupt the model’s decision boundaries. For DEEPTASTER, we have chosen to use FGSM for its simplicity and satisfactory performance in generating adversarial images. We employ Foolbox [34], a commonly used library, to facilitate our experiments.

FGSM is a gradient-based adversarial attack algorithm [12]. Let us consider x as the original image and ∇ as a slight perturbation applied to x , which leads to the creation of an adversarial sample, denoted as \bar{x} . The training process seeks to maximize the loss function $J(x, y)$ to derive the adversarial sample \bar{x} , which no longer belongs to class y . The entire optimization process has to fulfill the L_∞ constraint $\|\bar{x} - x\|_\infty \leq \epsilon$. Accordingly, FGSM adversarial examples can be produced using the following equation:

$$\bar{x} = x + \epsilon \cdot \text{sgn}(\nabla_x J(f(x), y)) \quad (1)$$

Here, sgn is the sign function, $J(f(x), y)$ is the loss function of the model’s prediction for input x and the true label y , and ϵ is a small constant which controls the magnitude of the perturbation.

2.2 Discrete Fourier Transform (DFT)

The Discrete Fourier Transform (DFT) is a tool used to transform a sequence of numbers $\{x_0, x_1, \dots, x_N\}$ in the time domain into another sequence of numbers $\{y_0, y_1, \dots, y_N\}$ in the frequency domain. The transformation is achieved through the equation $y_k = \sum_{n=0}^N x_n \cdot e^{-\frac{i2\pi}{N}kn}$. Upon application to an image, the DFT translates its spatial content into a frequency spectrum. This spectrum can be represented as an image, often revealing patterns imperceptible in the spatial domain.

Adversarial images are typically noisier than natural images. These noises are a byproduct of the adversarial perturbation process. We found that the frequency domain shows these noises more clearly than the spatial domain. This is aligned with the findings of Harder *et al.* [16], who showed that small changes in the spatial domain are hard to detect because they appear random. However, the same changes can lead to systematic changes in the frequency domain, which are then detectable. Our experimental results show that observing adversarial perturbations in the Fourier domain can be more advantageous than the spatial domain in identifying DNN models built with a specific dataset.

3 THREAT MODEL

Overview. We examine the issue of dataset leakage in the context of DNN models. High-profile incidents, such as the theft of over 5.6 million fingerprint records from the US Office of Personnel Management (OPM) [13] and the “Capital One” data breach [31],

have highlighted the risks posed by both external and insider threats. If a dataset is leaked, an adversary could use it to create a new DNN model or improve an existing one. In another scenario, an adversary could steal models from the victim’s private cloud or use the MLaaS API of the victim model to extract them. In this case, the adversary could use fine-tuning, pruning, and transfer learning to improve the performance of the stolen model and hide the signs of theft. Therefore, we considered the following threat models.

Assumptions. We make the following assumptions: (a) The adversary can steal the dataset used to train a victim’s DNN model or the model itself. Malicious insiders can access unencrypted databases, while external threats could employ advanced techniques (e.g., SQL injection) to obtain plaintext data from databases. (b) The adversary aims to build a model with the illegally obtained dataset and evade copyright verification. (c) The surrogate model, as crafted by the adversary, provides sufficient accuracy that the adversary can profit from its sale or commercialisation.

Settings. In our experiments, we considered the following different adversarial settings. Table 1 summarises these attacks along with the assumptions about the attacker’s access level. None of the existing DNN IP protection mechanisms has considered attack scenarios (1), (2), (7), and (8).

Table 1: Summary of Adversarial Settings. The “Access” column indicates the access required by an attacker.

N	Attack	Access	
		Dataset	Model
1	Multi-Architecture Attack (MAA)	Full	None
2	Data Augmentation Attack (DAA)	Full	None
3	Model Retraining Attack (SAA)	Partial	None
4	Transfer Learning Attack (TLA)	None	Full
5	Model Fine-tuning Attack (MFA)	Partial	Full
6	Model Pruning Attack (MPA)	Full	Full
7	Data Augmentation and Transfer Learning Attack (DATLA)	Full	Full
8	Transfer Learning with Pretrained mode Attack (TLPA)	Full	None

(1) **Multi-Architecture Attack (MAA).** The attacker steals the victim’s dataset. The attacker trains a model with an architecture different from the original victim model using the stolen data.

(2) **Data Augmentation Attack (DAA).** The attacker steals the victim’s dataset and creates a new one by combining the stolen data with new data. The attacker is aware of the structure of the victim’s model. The attacker trains a model that mirrors the victim’s model structure using the combined data.

(3) **Same Architecture Attack (SAA)** [27]. The attacker steals part of the victim’s dataset. The attacker is aware of the structure of the victim’s model. The attacker trains a model mirroring the victim’s model structure using the stolen data.

(4) **Transfer Learning Attack (TLA)** [27]. The attacker steals the victim’s model. The attacker uses transfer learning to fine-tune it on another dataset.

(5) **Model Fine-tuning Attack (MFA)** [27]. The attacker steals part of the victim’s dataset and the victim’s model. The attacker fine-tunes the model using the stolen dataset.

(6) **Model Pruning Attack (MPA)** [27]. The attacker steals the victim’s model. The attacker prunes the model and redistributes it.

(7) **Data Augmentation and Transfer Learning Attack (DATLA).** The attacker steals the victim’s dataset and model. They then create

a new dataset by combining the stolen data with new data. The attacker uses transfer learning to fine-tune the stolen model using the combined dataset.

(8) **Transfer Learning with Pretrained model Attack (TLPA).** The attacker steals the victim’s dataset. The attacker uses transfer learning to fine-tune the pretrained model with another dataset on the victim’s dataset.

4 DEEPTASTER SYSTEM DESIGN

In this section, we present DEEPTASTER, a DNN IP tracking tool that verifies whether an attacker’s model has been trained using a victim’s dataset or model. We first present an overview of the system design and then detail the system’s three main components: adversarial perturbation generation and transformation, classifier generation, and verification.

4.1 DEEPTASTER Overview

DEEPTASTER operates in two stages: (a) constructing a classifier using adversarial images in the DFT domain from a set of representative models trained on the victim dataset, and (b) verifying a suspect model by testing it using its adversarial images in the DFT domain. Figure 1 provides a schematic representation of this process.

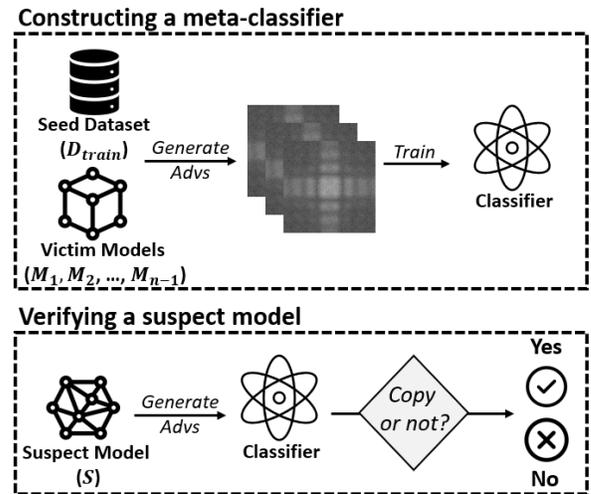


Figure 1: Overview of DEEPTASTER.

4.1.1 **Adversarial Perturbation Generation and Transformation.** Given a victim dataset, the dataset proprietor uses Algorithm 1 to generate adversarial images from a target model M constructed with the dataset. That is, the FGSM attack is performed on the target model M to generate adversarial images as dictated by Equation 1.

We use the Foolbox library [34] to generate adversarial images using the FGSM method. We use the ℓ_2 -norm metric to calculate the distance and set the epsilon value to 0.09. The same seed images are used for all victim models tested to generate adversarial images. When the seed image domain differs from the victim model’s image domain, we apply the FGSM method by re-labeling the seed images with the prediction value of the model. Adversarial images are

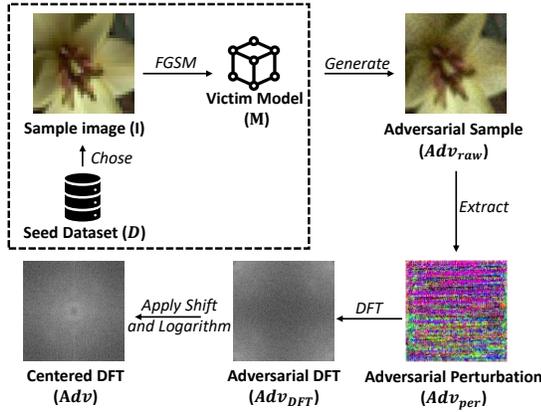
Algorithm 1 Adversarial Image Generation and Transformation.**Input:** Sample image I and target model M **Output:** Adversarial DFT image Adv

```

1: procedure GenerateAdv( $M, I$ )
2:    $Adv_{raw} \leftarrow FGSM(M, I)$ 
3:    $Adv_{per} \leftarrow Adv_{raw} - I$ 
4:    $Adv_{DFT} \leftarrow FourierTransform(Adv_{per})$ 
5:    $Adv \leftarrow Log(Shift(Adv_{DFT}))$ 
6:   return  $Adv$ 
7: end procedure

```

selected when they successfully fool a victim model into producing an incorrect prediction. The adversarial perturbation Adv_{per} is the pixel-wise difference between the original image I and its adversarial image Adv_{raw} . We then transform the adversarial image Adv_{raw} into the frequency domain, resulting in the adversarial DFT image Adv_{DFT} . To more accurately capture the characteristics of the dataset intelligence using the adversarial DFT images, we generate a *centered* DFT image, denoted as Adv , by applying a shift operation and a logarithm function sequentially to the adversarial DFT image Adv_{DFT} . Figure 2 illustrates the process of generating a centered DFT image. However, for improved readability, we will refer to Adv as the adversarial DFT image throughout the rest of the paper instead of using the term centered DFT images.

**Figure 2: Adversarial Image Generation.**

4.1.2 Classifier Construction. We have constructed a one-class classifier that is trained on adversarial DFT images generated from multiple model architectures, all of which are trained on the same dataset. Our goal is to build a robust detector that can determine whether a stolen dataset has been used to build a target model, even when the target model differs from the original one. To achieve this goal, we have chosen the Deep Support Vector Data Description (DeepSVDD) model [35] as our classifier from among several one-class classification models. DeepSVDD is an anomaly detection model that is widely used in different domains. It uses a deep neural network structure to find a hypersphere in the feature space. Therefore, we use DeepSVDD to extract data features from

adversarial discrete Fourier transform (DFT) images for models trained on a specific (victim) dataset. These features can then be used effectively by DeepSVDD to identify such models.

We construct a one-class classifier for DEEPTASTER using a seed dataset D and a set of victim models V , which are trained on the victim dataset. To determine the optimal threshold of the classifier, we use a set of benign models B_{val} trained on a different dataset than V . We split D into two parts: D_{train} and D_{val} . Similarly, we split V into V_{train} and V_{val} . We use D_{train} and V_{train} to train the one-class classifier and D_{val} , V_{val} , and B_{val} to optimize the classifier’s threshold τ .

To train the classifier, we generate adversarial samples Adv_{train} from D_{train} using Algorithm 1 for the victim models in V_{train} . We then use Adv_{train} to train a DeepSVDD model as a one-class classifier.

To validate the classifier, we generate adversarial samples Adv_{val} from D_{val} for the victim models in V_{val} using Algorithm 1. We then evaluate the classifier’s ability to distinguish between V_{val} and B_{val} using Adv_{val} to determine the threshold τ . The classifier is used to check whether a suspect model is built on V or another dataset, and we set the threshold τ to maximize the classifier’s accuracy.

For testing, the one-class classifier uses the threshold τ to verify if a suspect model is built on the victim’s dataset. The details of the model verification process are given in Section 4.1.3. Using a larger set of victim models with different architectures leads to higher accuracy for the one-class classifier.

We surmise that any random images could be used as seed images for the victim model. In our experiments, we randomly selected samples from CIFAR10 as seed images to generate adversarial images for models built on not only CIFAR10 but also other datasets such as MNIST and Tiny-ImageNet. However, analyzing the effects of seed samples and developing an algorithm to select effective seed samples would be an interesting area for future work.

4.1.3 Verification. We test whether a suspect model S was built using the victim dataset with the constructed one-class classifier using D_{val} . This dataset D_{val} is used both for optimizing the threshold and generating the adversarial DFT images Adv_{test} for verification. Algorithm 2 provides a detailed description of the verification procedure of DEEPTASTER. The algorithm inputs include the classifier \mathcal{M} , its threshold value τ , the dataset D_{val} , and a suspect model S . The output is a verification result indicating whether the suspect model contains part of the victim dataset.

In order to verify the suspect model S , we initially generate the adversarial DFT images from D_{val} for S using the steps described in Section 4.1.1, as shown in lines 1–3. Next, we put these images to the test by feeding them into the one-class classifier, as shown in lines 4–7. If the fraction of samples that fall below the classifier’s threshold τ exceeds 50%, we classify the suspect model as having been trained on stolen data, as shown in lines 8–12. The default criterion for DEEPTASTER is set at 50% to simplify the decision-making process and ensure an unbiased assessment. The “theft image detection rate” refers to the proportion of tested adversarial samples that fall below the classifier’s threshold τ .

5 EXPERIMENTS

We implemented DEEPTASTER as a self-contained toolkit in Python. In this section, we evaluate the performance of DEEPTASTER against

Algorithm 2 Validation using DEEPTASTER.

Input: Classifier \mathcal{M} , the threshold τ , the dataset D_{val} , and the suspect model S .

Output: Verification results

```

1:  $Adv_{test} \leftarrow GenerateAdv(S, D_{val})$ 
2:  $X \leftarrow 0$ 
3:  $k \leftarrow len(Adv_{test})$ 
4: while  $k \neq 0$  do
5:    $X \leftarrow X + \mathcal{M}(Adv_{test}[k]) \leq \tau$ 
6:    $k \leftarrow k - 1$ 
7: end while
8: if  $X > len(Adv_{test}) * \frac{1}{2}$  then
9:    $S$  is a stolen model
10: else
11:    $S$  is a benign model
12: end if

```

an extensive list of eight attacks mentioned in Section 3. Some of these attacks, such as fine-tuning and pruning, are well-studied in watermarking [27]. We also examine DEEPTASTER against more challenging adaptive attack scenarios such as transfer learning, re-training, and the most challenging – multi-architecture – which has never been considered before in the literature. To ensure the generalisability of DEEPTASTER, we generate three classifiers that track CIFAR10, MNIST, and Tiny-ImageNet. We also compare our results to the best state-of-the-art fingerprinting technique, DEEPTJUDGE [8].

5.1 Experimental Setup

Datasets and victim models. We use four datasets, including CIFAR10 [21], MNIST [25], Tiny-ImageNet [23], and ImageNet [10]. The first three datasets are used as both victim and suspect datasets. The ImageNet dataset is used as a suspect dataset only. All datasets are image classification datasets with a varying number of classes, ranging from 10 classes in CIFAR10 and MNIST to up to 1000 in ImageNet, as described in Appendix B. We note that we use only half of the Tiny-ImageNet dataset (i.e., 100 classes) to run the experiments in order to reduce the experimental computation time.

We trained victim models using three popularly used DNN architectures, ResNet18 [17], VGG16 [36], and DenseNet161 [18], on each of the victim datasets. Appendix C provides a summary of the information about those models.

5.2 Classifier Evaluation Settings

Training configuration. We utilized the procedures described in Section 4.1.2 to construct a classifier capable of detecting whether a proprietary dataset has been used to build a model. In all experiments, we used a seed dataset consisting of 1888 randomly chosen images from CIFAR10, regardless of the datasets utilized for creating victim and suspect models. Subsequently, we generated 1888 adversarial DFT images for each victim model. These adversarial DFT images were split into two groups: 1600 images were used to train the classifier, and the remaining 288 images were used to determine the classifier’s threshold and for testing purposes. We discuss the effects of training dataset size and dimensions in Appendix D.

Effects of the threshold for the classifier. In Section 4.1.2, we discussed the process for determining the optimum threshold for a classifier using the validation dataset. We conducted experiments to investigate how the threshold value of the classifier impacts the performance of DEEPTASTER using nine models constructed from three datasets and three model architectures. We built a dedicated classifier for each dataset and evaluated its performance using balanced accuracy. Appendix E shows the relationship between each classifier’s threshold value and balanced accuracy. We observed that increasing the threshold value for all classifiers tends to lead to a significant increase in balanced accuracy, reaching a certain point (i.e., 0.0025 for CIFAR10, 0.0005 for MNIST, and 0.003 for Tiny-ImageNet) before decreasing. We selected the threshold value for each classifier that maximises its balanced accuracy.

5.3 Resilience against Data/Model Theft Attacks

We evaluate the resilience of DEEPTASTER against the eight attack scenarios presented in Section 3. We repeated each attack scenario ten times. The mean values for Model Accuracy (Model Acc.) and Theft image detection Rate (Theft Image Rate) are presented, along with the standard deviation values in parentheses. The same format is used in the remaining tables throughout this paper.

5.3.1 Multi-Architecture Attack (MAA). We test DEEPTASTER against the MAA scenario. In this scenario, the attacker uses part of the stolen dataset to train a model with an architecture different from the original victim model. We chose three victim datasets – CIFAR10, MNIST, and Tiny-ImageNet – and trained each of them using three different model architectures (ResNet18, VGG16, and DenseNet161). We designate one dataset as the victim and assign the other datasets as benign. Additionally, the models trained on ImageNet [10] are used as benign models for MNIST and CIFAR10. **Efficacy.** Table 2 presents DEEPTASTER’s performance against the MAA scenario. DEEPTASTER completely distinguishes attack cases from benign cases for all datasets. The results demonstrate that DEEPTASTER is highly effective against MAA under both stolen and benign scenarios. Figure 3 visually demonstrates the effectiveness of a classifier used for DEEPTASTER. The figure displays the distribution of output scores produced by the classifier for CIFAR10 across 12 distinct models, each representing a combination of the three architectures (ResNet18, VGG16, and DenseNet161) and the four datasets (CIFAR10, MNIST, Tiny-ImageNet, and ImageNet). The bold line in the figure represents the chosen threshold for the classifier, which is trained on the CIFAR10. The classifier can successfully distinguish the DNN model using the CIFAR10 dataset from the models using the other datasets.

Remark 1: DEEPTASTER is effective in detecting DNN model theft attacks across various model architectures.

5.3.2 Data Augmentation Attack (DAA). In the DAA scenario, we assume that an attacker tries to obscure the use of a stolen dataset by augmenting it with new data. For example, if the CIFAR10 dataset is used as a victim dataset, we added five new classes from the CIFAR100 dataset to create a new dataset dubbed CIFAR15. The added classes include apples, bicycle, can, roses, and clock, which are distinct from any classes in the original CIFAR10 dataset. We

Table 2: MAA results for CIFAR10, MNIST, Tiny-ImageNet classifiers with 12 suspect models, each representing a combination of three architectures – ResNet18 (RN), VGG16 (VGG), and DenseNet161 (DN) and four datasets – CIFAR10, MNIST, Tiny-ImageNet, and ImageNet. The “Copy Detection (%)” field values represent the successful copy detection rate. The same format is used in the remaining tables throughout the paper.

Victim	Suspect	Ground Truth	Architecture	Theft Image Rate	Copy Detection (%)
CIFAR10	CIFAR10	Stolen	RN	93.01 (5.34)	100
			VGG	84.65 (9.72)	100
			DN	94.55 (4.11)	100
	MNIST	Benign	RN	9.17 (16.02)	100
			VGG	0.0 (0.0)	100
			DN	9.8 (11.48)	100
	Tiny-ImageNet	Benign	RN	7.15 (8.01)	100
			VGG	0.21 (0.52)	100
			DN	3.64 (4.33)	100
	ImageNet	Benign	RN	4.24 (7.25)	100
			VGG	0.14 (0.23)	100
			DN	6.88 (12.94)	100
MNIST	MNIST	Stolen	RN	98.13 (3.51)	100
			VGG	92.40 (9.62)	100
			DN	95.97 (7.13)	100
	CIFAR10	Benign	RN	4.69 (10.66)	100
			VGG	0.0 (0.0)	100
			DN	0.28 (0.72)	100
	Tiny-ImageNet	Benign	RN	1.35 (2.14)	100
			VGG	0.0 (0.0)	100
			DN	0.14 (0.42)	100
	ImageNet	Benign	RN	5.14 (9.02)	100
			VGG	10.03 (15.90)	100
			DN	6.27 (13.26)	100
Tiny-ImageNet	Tiny-ImageNet	Stolen	RN	92.71 (9.04)	100
			VGG	72.61 (16.96)	100
			DN	90.70 (7.07)	100
	CIFAR10	Benign	RN	6.01 (13.51)	100
			VGG	11.35 (13.15)	100
			DN	4.20 (7.21)	100
	MNIST	Benign	RN	1.32 (2.61)	100
			VGG	1.36 (3.52)	100
			DN	7.43 (12.62)	100

train ResNet18 using the augmented CIFAR15 dataset as the attack case. The goal here is to evaluate whether DEEPTASTER can detect the usage of stolen CIFAR10 data, even when this data has been augmented with new classes. To test the effectiveness of DEEPTASTER at different stages of model training, we examine the performance of DEEPTASTER at five different epochs (20, 40, 60, 80, and 100), which can reveal if DEEPTASTER’s detection ability changes as the model becomes more trained.

Efficacy. Table 3 presents DEEPTASTER’s performance against the DAA scenario. The theft image detection rate slightly increases as the training epoch of the model increases. While DEEPTASTER failed to detect one or two attack cases, DEEPTASTER is still effective against DAA for large training epochs.

Remark 2: DEEPTASTER is effective in detecting cases where a model has been trained on a mixed dataset of stolen data and new data, especially when using large training epochs.

5.3.3 Same Architecture Attack (SAA). In the SAA scenario, an attacker trains the ResNet18 model on 10%, 30%, 50%, 70%, and 90%

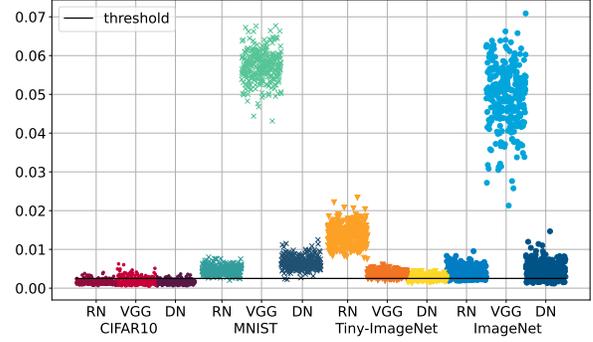


Figure 3: Distribution of output scores produced by the classifier for CIFAR10 across 12 different models, each representing a combination of three architectures – ResNet18 (RN), VGG16 (VGG), and DenseNet161 (DN) and four datasets – CIFAR10, MNIST, Tiny-ImageNet, and ImageNet. The bold line represents the threshold chosen for DEEPTASTER.

Table 3: DAA results for CIFAR10 classifier.

Suspect	Ground Truth	Epochs	Model Acc.	Theft Image Rate	Copy Detection (%)
CIFAR10	Stolen	20	72.61 (0.39)	77.12 (35.66)	80
		40	72.56 (0.69)	74.10 (33.34)	80
		60	72.67 (0.59)	72.01 (33.34)	90
		80	72.5 (0.57)	86.88 (16.50)	90
		100	72.67 (0.49)	77.99 (26.51)	90

of a victim dataset. We use CIFAR10 and MNIST, respectively, as victim datasets. We split the dataset uniformly, including an equal number of samples from every class. We examine the performance of DEEPTASTER at 10 epochs.

Efficacy. Table 4 presents the performance of DEEPTASTER against the SAA scenario. The experimental results indicate that for the MNIST dataset, DEEPTASTER can perfectly detect attacks when the portion of the stolen dataset used for training the model is 30% or more. However, in the case of CIFAR10, only 80% of the attacks were detected when the portion of the stolen dataset used for training was 70%. This suggests that the likelihood of an attack being detected by DEEPTASTER increases with the amount of stolen data used for training. This observation is consistent with our intuition, as we expect a model to exhibit more characteristics of the stolen data when a larger portion is used for training.

Remark 3: DEEPTASTER is effective in detecting DNN model theft attacks when a significant portion (more than 70%) of the stolen dataset is used for training the model.

5.3.4 Transfer Learning Attack (TLA). In the TLA scenario, an attacker steals the victim’s ResNet18 model trained on CIFAR10 and performs transfer learning with the MNIST dataset using a learning rate of 0.1 while freezing the lower 30% layer. We examine the performance of DEEPTASTER at four different epochs (2, 4, 6, and 8). **Efficacy.** Table 5 presents the performance of DEEPTASTER against the TLA scenario. DEEPTASTER failed to detect four attack cases

Table 4: SAA results for CIFAR10 and MNIST classifier when the different amount of dataset is used.

Victim	Suspect	Ground Truth	Used Dataset Size (%)	Theft Image Rate	Copy Detection (%)
CIFAR10	CIFAR10	Stolen	10	48.86 (24.56)	40
			30	57.05 (27.07)	70
			50	64.41 (29.02)	60
			70	80.03 (23.32)	80
			90	79.10 (26.55)	90
MNIST	MNIST	Stolen	10	65.56 (11.26)	90
			30	80.35 (9.53)	100
			50	82.40 (8.59)	100
			70	93.61 (5.41)	100
			90	95.14 (3.82)	100

when the attacker trained the model for only two epochs. However, DEEPTASTER was able to completely distinguish attack cases as stolen at 8 epochs.

Table 5: TLA results for CIFAR10 classifier.

Suspect	Ground Truth	Epochs	Model Acc.	Theft Image Rate	Copy Detection (%)
CIFAR10 to MNIST	Stolen	2	98.42 (0.21)	54.31 (11.10)	60
		4	98.65 (0.12)	62.09 (6.56)	100
		6	98.50 (0.27)	64.76 (12.08)	90
		8	98.92 (0.21)	69.93 (7.72)	100

Remark 4: DEEPTASTER shows robust effectiveness in identifying cases where a stolen model has been further trained with new data through transfer learning when the model is fully trained.

Table 6: MFA results for CIFAR10 classifier.

Suspect	Ground Truth	Dataset Size	Model Acc.	Theft Image Rate	Copy Detection (%)
CIFAR10	Stolen	500 (0.01%)	74.67 (0.39)	99.93 (0.14)	100
		1000 (0.02%)	74.77 (0.22)	99.93 (0.14)	100
		2500 (0.05%)	75.06 (0.04)	99.83 (0.23)	100
MNIST	Benign	500 (0.01%)	99.46 (0.0)	0.0 (0.0)	100
		1000 (0.02%)	99.48 (0.0)	0.0 (0.0)	100
		2500 (0.05%)	99.48 (0.0)	0.0 (0.0)	100

5.3.5 Model Fine-tuning Attack (MFA). In the MFA scenario, an attacker steals the ResNet18 model trained on CIFAR10 and the CIFAR10 dataset itself. The attacker then uses fine-tuning to train the stolen model using part of the CIFAR10 dataset, consisting of either 500, 1000, or 2500 samples, respectively, using the learning rate of 0.00005. For comparison, a benign model is trained on only the MNIST dataset using the same architecture. We examine the performance of DEEPTASTER at 60 epochs.

Efficacy. Table 6 presents the performance of DEEPTASTER against the MFA scenario. DEEPTASTER completely distinguishes all MFA cases from benign cases, irrespective of the used dataset size.

Remark 5: DEEPTASTER shows robust effectiveness in identifying cases where a stolen model has been further trained using a subset of the dataset used to train the stolen model.

5.3.6 Model Pruning Attack (MPA). In the MPA scenario, an attacker steals the ResNet18 model trained on CIFAR10 and prunes 20%, 40%, and 60% of the stolen model. Then the attacker fine-tunes the model for five epochs with a learning rate of 0.00005. For CIFAR10, we examine the performance of DEEPTASTER at 5 epochs. For comparison, a benign model is trained on only the MNIST dataset using the same architecture. For MNIST, we examine the performance of DEEPTASTER at 5 epochs.

Efficacy. Table 7 presents the performance of DEEPTASTER against the MPA scenario. DEEPTASTER completely distinguishes all MPA instances from benign cases, regardless of the degree of pruning applied to the model.

Table 7: MPA results for CIFAR10 classifier with three positive models and three negative models.

Suspect	Ground Truth	Pruned Percentage	Model Acc.	Theft Image Rate	Copy Detection (%)
CIFAR10	Stolen	20	63.09 (0.12)	100 (0.0)	100
		40	42.73 (0.07)	99.97 (0.11)	100
		60	21.44 (0.08)	100 (0.0)	100
MNIST	Benign	20	98.53 (0.01)	4.62 (7.22)	100
		40	87.54 (0.11)	1.25 (3.75)	100
		60	43.18 (0.13)	4.83 (14.48)	100

Remark 6: DEEPTASTER is robust against MPA regardless of the percentage of neurons pruned.

5.3.7 Data Augmentation and Transfer Learning Attack (DATLA). In the DATLA scenario, an attacker creates the CIFAR15 dataset using the method explained in Section 5.3.2. The attacker then uses a stolen ResNet18 model trained on the CIFAR10 dataset and fine-tunes it on the CIFAR15 dataset the attacker created. We examine the performance of DEEPTASTER at five different epochs (20, 40, 60, 80, and 100).

Efficacy. Table 8 shows DEEPTASTER’s performance against the DATLA scenario. DEEPTASTER failed to distinguish one or three cases when the training epoch is 20, 40, 60, and 80. However, DEEPTASTER successfully distinguished all attack cases when the training epoch is 100. This highlights the effectiveness of DEEPTASTER when using larger training epochs while also suggesting areas for potential improvement in handling cases with minimal training epochs.

Remark 7: DEEPTASTER is effective in detecting cases where a stolen model has been trained on a mixed dataset that combines stolen data with new data while suggesting areas for potential improvement in handling cases with minimal training epochs.

5.3.8 Transfer Learning with Pretrained model Attack (TLPA). In the TLPA scenario, an attacker uses a pretrained model on a public dataset, such as MNIST. The attacker then steals the CIFAR10 dataset of 60,000 samples and uses transfer learning to fine-tune the pretrained model using the stolen dataset. Similar to the TLA scenario, the attacker performs transfer learning using a learning rate of 0.1 while freezing the lower 30% layers. We examine the performance of DEEPTASTER at 10 epochs by varying the dataset size used to generate the pretrained model.

Table 8: DATLA results for CIFAR10 classifier.

Suspect	Ground Truth	Epochs	Model Acc.	Theft Image Rate	Copy Detection (%)
CIFAR10	Stolen	20	73.51 (0.74)	84.62 (27.83)	90
		40	73.53 (0.50)	83.26 (24.61)	90
		60	73.56 (0.47)	82.92 (19.55)	90
		80	73.63 (0.51)	62.60 (36.40)	70
		100	73.24 (0.51)	81.25 (11.36)	100

Efficacy. Table 9 shows DEEPTASTER’s performance against the TLPA scenario. The theft Imagetection rate decreases as the dataset size used for the pre-trained model increases. DEEPTASTER can completely detect TLPA attacks when at least 30% of the MNIST dataset is used for pretraining. However, if the dataset used for the pretrained model exceeds 50%, DEEPTASTER may fail to detect some stolen cases.

Table 9: TLPA results for CIFAR10 classifier.

Suspect	Ground Truth	Used Dataset Size (%)	Model Acc.	Theft Image Rate	Copy Detection (%)
MNIST to CIFAR10	Stolen	10	67.66 (1.13)	99.51 (1.46)	100
		30	68.43 (0.63)	99.58 (0.62)	100
		50	68.74 (0.86)	77.99 (23.71)	90
		70	68.74 (0.86)	84.65 (13.94)	100
		90	69.22 (0.64)	85.28 (23.90)	80

Remark 8: DEEPTASTER shows robust effectiveness in identifying cases where a pretrained model has been further trained using a stolen dataset via transfer learning.

5.4 Comparison with Existing Fingerprinting

Comparison settings. We conduct experimental comparisons with DEEPJUDGE [8], the state-of-the-art fingerprinting technique. DEEPJUDGE generates four metrics for white-box evaluation and two metrics for black-box evaluation. It uses majority voting, where 3 out of 4 metrics have to produce values $<$ threshold to support the correct final judgment of being stolen. DEEPJUDGE has been designed to provide architecture-dependent protection; namely, all model parameters need to be the same, including the number of classes. On the other hand, DEEPTASTER is designed to be architecture insensitive to enable the dataset intelligence to be tracked even when the model architecture is changed. Hence, DEEPJUDGE is not able to detect MAA as stolen models owing to its design limitation. Therefore, we compare DEEPJUDGE with DEEPTASTER in the other seven attacks in addition to direct cloning. For the transfer learning attack, we observe that DEEPJUDGE only considered the same number of classes between the original model and the transfer-learned model, which might not always be the case. Thus, we use the scripts released with DEEPJUDGE and apply small modifications to run the data augmentation and transfer learning attacks when the number of classes is different. We set the target model of DEEPJUDGE as the ResNet18 model trained on CIFAR10, MNIST, and Tiny-ImageNet. For the CIFAR10 dataset, we test with five attack models, using a ResNet18 model trained on MNIST and Tiny-ImageNet as the “Benign” cases. In other cases, we only check the theft image detection rate for two negative cases trained on other datasets. We used the

optimal threshold for DeepJudge, which was optimized based on the testing results. These thresholds may be infeasible to achieve in practice, but it is ideal for comparison purposes. We only used four white-box metrics from DeepJudge.

Table 10: Detection results of DEEPJUDGE [8]. DEEPJUDGE uses a majority voting mechanism to determine whether a stolen model is being used. Specifically, DEEPJUDGE classifies a model as stolen if at least three out of four metrics produce values $<$ threshold.

Victim	Ground Truth	Suspect	Metric1	Metric2	Metric3	Metric4	Copy Detection (%)	
CIFAR10	Threshold		0.0137	0.0542	0.3544	0.4255		
	Stolen	CIFAR10	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	100	
		CIFAR10 DAA	0.0055 (0.0032)	0.0241 (0.0116)	0.5681 (0.3179)	0.5771 (0.3503)	40	
		CIFAR10 SAA	0.0050 (0.0023)	0.0213 (0.0084)	1.0198 (0.3477)	1.0123 (0.3324)	10	
		CIFAR10 TLA	0.0306 (0.0)	0.1215 (0.0002)	0.3463 (0.0)	0.3188 (0.0)	0	
		CIFAR10 MFA	0.0005 (0.0)	0.0018 (0.0005)	0.0 (0.0)	0.0 (0.0)	100	
		CIFAR10 MPA	0.0108 (0.0001)	0.0433 (0.0005)	0.0830 (0.0)	0.0775 (0.0)	100	
		CIFAR10 DATLA	0.0071 (0.0027)	0.0293 (0.010)	0.1958 (0.0717)	0.1846 (0.0687)	100	
		CIFAR10 TLPA	0.0101 (0.0022)	0.0453 (0.0076)	1.5621 (0.4715)	1.5654 (0.4720)	0	
		Benign	MNIST	0.0925 (0.0331)	0.3724 (0.1317)	0.8535 (0.3554)	0.8597 (0.3503)	100
			MNIST SAA	0.0806 (0.0328)	0.2399 (0.1316)	1.1233 (0.5622)	1.1528 (0.5460)	100
	MNIST MFA		0.0846 (0.0)	0.3408 (0.0)	1.003 (0.0)	1.032 (0.0)	100	
	MNIST MPA		0.664 (0.0)	0.2705 (0.0004)	0.8244 (0.0)	0.8810 (0.0)	100	
	Tiny ImageNet		0.0164 (0.0050)	0.0643 (0.0193)	1.2067 (0.4641)	1.2083 (0.4415)	90	
	Threshold		0.0407	0.1846	0.3022	0.3094		
	MNIST	Stolen	MNIST	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	100
			MNIST SAA	0.042 (0.0279)	0.17 (0.1106)	0.489 (0.3336)	0.491 (0.3338)	40
			CIFAR10	0.0964 (0.0343)	0.3860 (0.1371)	0.8980 (0.5226)	0.9060 ((0.5413))	100
		Benign	Tiny ImageNet	0.1027 (0.3327)	0.4111 (0.1329)	0.9823 (0.2940)	0.9848 (0.2979)	100
			Threshold		0.0020	0.0095	0.3660	0.3660
Tiny-ImageNet	Stolen	Tiny ImageNet	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	100	
		Tiny ImageNet SAA	0.0030 (0.0024)	0.0131 (0.0091)	0.5111 (0.3003)	0.5111 (0.3003)	30	
	Benign	CIFAR10	0.0052 (0.0019)	0.0223 (0.0069)	0.9156 (0.3672)	0.9156 (0.3672)	100	
		MNIST	0.0047 (0.0013)	0.0196 (0.0052)	0.8433 (0.3290)	0.8433 (0.3290)	100	

Results. Table 10 presents the comparison results. For CIFAR, DEEPJUDGE was ineffective in identifying most instances of theft against DAA (40%), SAA (10%), TLA (0%), and TLPA (0%). However, under the same settings, DEEPTASTER effectively detected those attacks (DAA (90%), SAA (90%), TLA (100%), and TLPA (100%)). For MNIST, the detection rate dropped to 40% under the SAA scenario. Similarly, for Tiny-ImageNet, the detection rate dropped to 30% under the SAA scenario.

5.5 Robustness of DEEPTASTER with a Large-Sized Model Architecture

To evaluate the robustness of DEEPTASTER on large-sized model architectures, we used AlexNet [22] as a representative architecture and Tiny-ImageNet as the victim dataset. We constructed a classifier with four victim models: AlexNet, ResNet18, VGG16, and DenseNet161, all trained on the Tiny-ImageNet dataset. We then

evaluated the constructed classifier using twelve models, including the four victim and eight benign models trained on CIFAR10 or MNIST. As a result, DEEPTASTER remains effective in detecting adversarial images for large-scale model architectures such as AlexNet (see Appendix F). We only conducted this experiment once due to the significant time required for experiments on AlexNet.

5.6 Ablation Study on DFT

We conducted an ablation study to assess DFT’s impact on DEEPTASTER’s performance. We built a classifier identical to DEEPTASTER but without DFT, and tested its performance under the Multi-Architecture Attack scenario described in Section 5.3.1. We used four datasets (CIFAR10, MNIST, Tiny-ImageNet, and ImageNet) and three architectures (ResNet18, VGG16, and DenseNet161) for this evaluation. Table 11 indicates a significant degradation in accuracy when DFT is not applied. Without DFT, the classifier incorrectly classified all stolen cases as benign, except for one case in each CIFAR10 and MNIST dataset when the DenseNet161 architecture was used. Conversely, as shown in Table 2, DEEPTASTER could correctly detect benign and stolen cases with DFT. These findings underscore the critical role of DFT in maintaining the efficacy of DEEPTASTER.

Table 11: Ablation study on DFT results.

Victim	Suspect	Ground Truth	Architecture	Theft Image Rate	Copy Detection (%)
CIFAR10	CIFAR10	Stolen	RN	4.62 (11.82)	0
			VGG	5.35 (13.98)	0
			DN	11.25 (22.62)	10
	MNIST	Benign	RN	4.55 (10.19)	100
			VGG	0.0 (0.0)	100
			DN	1.39 (4.17)	100
	Tiny-ImageNet	Benign	RN	0.0 (0.0)	100
			VGG	0.14 (0.42)	100
			DN	1.04 (2.09)	100
	ImageNet	Benign	RN	15.97 (25.64)	90
			VGG	8.09 (17.51)	90
			DN	12.78 (23.02)	90
MNIST	MNIST	Stolen	RN	4.51 (10.05)	0
			VGG	0.0 (0.0)	0
			DN	11.08 (29.77)	10
	CIFAR10	Benign	RN	0.42 (0.95)	100
			VGG	0.17 (0.42)	100
			DN	2.40 (5.99)	100
	Tiny-ImageNet	Benign	RN	0.035 (0.105)	100
			VGG	0.0 (0.0)	100
			DN	0.35 (0.93)	100
	ImageNet	Benign	RN	1.15 (2.64)	100
			VGG	0.35 (0.71)	100
			DN	0.56 (1.35)	100
Tiny-ImageNet	Tiny-ImageNet	Stolen	RN	0.07 (0.14)	0
			VGG	0.0 (0.0)	0
			DN	0.10 (0.22)	0
	CIFAR10	Benign	RN	0.07 (0.14)	100
			VGG	0.04 (0.11)	100
			DN	0.04 (0.11)	100
	MNIST	Benign	RN	0.17 (0.32)	100
			VGG	0.35 (0.11)	100
			DN	0.0 (0.0)	100

5.7 Detection Latency Results

We assessed the detection latency of DEEPTASTER by analyzing the time taken by the classifier to identify a suspect model. The process consisted of three steps, as shown in Table 12.

Table 12: Mean time taken for each step in DEEPTASTER.

Step	Task	Time (Sec)
1	Adversarial DFT Generation	668.0
2	Classifier Training	766.8
3	Suspect Model Verification	9.78

The first step involved the generation of adversarial images in the DFT domain. This took approximately 668 seconds, equivalent to 0.3 seconds per image. The second step comprised training the classifier, taking about 766.8 seconds. It is important to note that the first two steps are one-off tasks, meaning they are not repeated for each verification. The final step was verifying the suspect model, taking around 9.78 seconds, equivalent to 0.03 seconds per image. The entire process, including generation, training, and verification, took approximately 1444.51 seconds. These latency results are comparable to those of DEEPTJUDGE, which took 1937.79 seconds.

We plan to optimize the verification process by reducing the number of adversarial images tested. We aim to use only a few images for verification. This approach should enable DEEPTASTER to verify a suspect model within 1 second while still maintaining reasonable detection accuracy.

6 DISCUSSION

Usefulness of DEEPTASTER. DEEPTASTER exhibits usefulness in various aspects. First, to ensure robustness, DEEPTASTER effectively detected DNN model theft attacks under various conditions. For example, DEEPTASTER was able to detect attacks on models trained on different datasets, attacked with different adversarial attack methods, and even attacked by models trained with a small subset of the victim dataset. Second, to ensure fidelity, DEEPTASTER has zero impact on model accuracy as it uses a fingerprinting technique rather than invasive watermarking. This means that DEEPTASTER can be used to protect models without sacrificing their performance. Third, to ensure efficacy, DEEPTASTER can distinguish most of the attack cases against benign models, except for some extreme cases (i.e., trained with small epochs or trained with less than 70% of the victim dataset). This means that DEEPTASTER is a highly effective tool for detecting DNN model theft attacks. Overall, DEEPTASTER is a robust, reliable, and effective tool for protecting machine learning models from theft.

Robustness against unseen architectures. DEEPTASTER was specifically designed to detect DNN model theft attacks across various architectures. However, our experiments suggest that DEEPTASTER may not be effective in detecting models trained using completely new or unseen architectures (see Table 13).

To evaluate the effectiveness of DEEPTASTER against models with an unseen architecture, we built a classifier using the CIFAR10 dataset as the victim dataset. We then tested the classifier’s robustness against six different models, three of which were trained using ResNet101 and the other three using SqueezeNet [19]. The evaluation results are presented in Table 13.

Unfortunately, our experiments revealed that DEEPTASTER was ineffective in detecting the use of a stolen dataset (CIFAR10) for both ResNet101 and SqueezeNet models. This is because the generated adversarial images are influenced not only by the dataset used but

Table 13: Performance of DEEPTASTER with unseen architectures when the victim dataset is CIFAR10.

Architecture	dataset	ground truth	Theft Image Rate	Copy Detection(%)
ResNet101	CIFA10	stolen	0.0	0
	MNIST	benign	0.0	100
	Tiny-ImageNet	benign	0.0	100
SqueezeNet	CIFA10	stolen	0.0	0
	MNIST	benign	0.0	100
	Tiny-ImageNet	benign	0.0	100

also by the underlying model architecture. Based on these findings, we conclude that DEEPTASTER may not be effective in detecting models trained using completely new or unseen architectures.

To overcome this limitation, we need to build a more diverse and comprehensive classifier that includes a broader range of architectures. As demonstrated in Section 5.5, our experimental results indicate that DEEPTASTER effectively detects dataset theft using AlexNet when the classifier is constructed using the same architecture (i.e., AlexNet). However, the effectiveness of DEEPTASTER declines when the attacker’s model is not included in the classifier construction phase. Nonetheless, it is worth noting that DEEPTASTER is currently the only fingerprinting scheme that can operate across multiple architectures.

Robustness against watermark removal attacks. Various watermark removal attacks have been introduced to evaluate the robustness of DNN watermarking techniques [3]. Traditional watermark removal attacks are ineffective because DeepTaster does not rely on specialized watermark-triggered input samples. However, a recent study by Guo *et al.* [15] introduced a new watermark removal attack technique that uses preprocessing. They discovered that watermark samples are less robust than normal samples and designed a preprocessing function to compromise the watermark verification output without affecting the normal output. This approach employs a series of transformations, such as scaling, embedding random imperceptible patterns, and spatial-level transformations, to effectively disable watermark-triggered input samples while maintaining the model’s accuracy. The preprocessing methods used in this scheme could make it more challenging to find effective adversarial examples for DeepTaster. We plan to explore this possibility in future work.

Robustness against adversarial training. To evaluate the impact of adversarial training on DEEPTASTER’s effectiveness, we conducted experiments using the ResNet18 architecture and the CIFAR10 dataset. We used an adversarial training method [4] than the method introduced by Madry et al. [30]. While Madry et al.’s method was originally designed to improve robustness, the method we used focuses on maintaining model accuracy when training on adversarial examples. As a result, our model maintains its accuracy after adversarial training. In this attack scenario, an attacker steals the CIFAR10 dataset and generates adversarial examples using the FGSM attack with an epsilon of 0.001. Our method, on the other hand, aims to maximize model accuracy. The attacker then trains the ResNet18 model on these adversarial examples, which comprise 1%, 3%, 5%, 7%, and 9% of the original CIFAR10 dataset, respectively. We assessed DEEPTASTER’s performance over 10 epochs, repeating the same experiments 10 times for each configuration to avoid bias.

Table 14 presents the results. We found that DEEPTASTER’s performance decreased as the proportion of adversarial examples in the training data increased. Even when only 1% of adversarial examples were used, DEEPTASTER failed once out of 10 attempts. When 5% or more of the training data consisted of adversarial examples, DEEPTASTER failed as many as four times. This result is not surprising, as DEEPTASTER relies on the characteristics of adversarial examples to determine the model’s decision boundary. When new adversarial examples are introduced during training, the existing examples used for fingerprinting struggle to accurately reflect the model’s decision boundary. To overcome this limitation, future research could explore pretraining victim models with adversarial training before applying DEEPTASTER.

Table 14: Performance of DEEPTASTER against adversarial training.

Suspect	Ground Truth	Used Dataset Size (%)	Model Acc.	Theft Image Rate	Copy Detection (%)
CIFAR10	Stolen	1	70.45 (0.96)	75.04 (16.70)	90
		3	70.15 (1.02)	66.74 (30.32)	80
		5	70.52 (0.88)	56.67 (26.60)	60
		7	70.53 (0.59)	54.83 (23.21)	60
		9	70.24 (0.64)	43.20 (31.46)	60

Comparison to membership inference attacks. DEEPTASTER and membership inference attacks are both methods to infer whether a target model contains samples from a dataset. Therefore, conventional membership inference attacks could potentially be used as a model fingerprinting technique.

To discuss the effectiveness of DEEPTASTER in detecting stolen datasets against membership inference attacks, we conducted experiments using TrajectoryMIA [26], a state-of-the-art membership inference attack method. TrajectoryMIA generates k different distillation models based on a target model and its shadow model, both trained on different datasets, during k epochs. The trajectory losses of a sample through each distillation model are computed to determine whether the sample is a member of the target model. TrajectoryMIA’s membership inference results can be used for model fingerprinting. If a sample from a victim dataset is a member of a suspect model, we can conclude that the suspect model was built on the victim dataset. To evaluate the performance of TrajectoryMIA, we used a VGG16 model trained on CIFAR10 as the suspect model. We used the same VGG16 model for its shadow and distillation models with the settings presented in [26]. The accuracy achieved by TrajectoryMIA on all CIFAR10 samples used as testing data was 60.52%. To compare the effectiveness of DEEPTASTER with TrajectoryMIA, we built a classifier of DEEPTASTER consisting of three architectures: ResNet18, VGG16, and DenseNet161. We evaluated the performance of DEEPTASTER using 288 adversarial images as testing samples for the victim model (i.e., VGG16 trained on CIFAR10), and it achieved a 98.26% accuracy. Although this comparison may not be entirely fair, we argue that utilizing a membership inference method in a straightforward manner might be insufficient for model fingerprinting. Moreover, DEEPTASTER is expected to be more efficient than TrajectoryMIA when testing multiple suspect models simultaneously since TrajectoryMIA requires creating several distillation models for each suspect model, whereas DEEPTASTER requires building a classifier only once for all suspect models.

Effects of the type of adversarial examples. The type of adversarial example can affect the performance of DEEPTASTER. To investigate this, we experimented with two popular attacks: PGD [29] and FGSM [12]. We compared the accuracy of the classifier on CIFAR10 using suspect models trained on CIFAR10, MNIST, and Tiny-ImageNet. DEEPTASTER using FGSM was about 7.73% more accurate than DEEPTASTER using PGD. Therefore, we recommend using the FGSM attack to build DEEPTASTER. Finding the most effective adversarial DFT images for a given suspect model is a challenging problem. As part of our future work, we plan to develop an algorithm to identify effective adversarial examples for DEEPTASTER.

On the high standard deviation of theft image detection rate: We observed a high standard deviation in the theft image detection rate, particularly for models trained from scratch compared to pre-trained models. This variance was prominent in specific attack scenarios: DAA (Table 3), SAA (Table 4), TLPA (Table 9), and adversarial training (Table 14). An exception is DATLA (Table 8), which requires transfer learning on a dataset with different classes, demanding more parameter adjustments. This may explain why DEEPTASTER can occasionally fail to detect theft images in these scenarios, resulting in an inflated standard deviation. However, the standard deviation is not as elevated when considering only successful instances of DEEPTASTER. For example, DAA reduced from 28.05% to 13.12%, SAA from 26.10% to 12.65%, DATLA from 23.95% to 11.5%, TLPA from 12.73% to 7.32%, and adversarial training from 25.66% to 12.96%. The high standard deviation can be attributed to the small sample size (10 experiments) and the inclusion of significantly lower detection rate outliers from failed cases. These findings indicate that the performance of DEEPTASTER can be inconsistent in specific attack scenarios. To address this limitation in the future, we may consider adopting more customized and optimized detection decision criteria than the fixed 50% threshold.

7 RELATED WORK

DNN watermarking. The first stream of related work uses watermarking to protect the copyright of DNN models [2, 9, 20, 24, 41, 47]. As in classical multimedia watermarking, DNN watermarking includes two stages: *embedding* and *verification*. In the *embedding* stage, the DNN model owner inserts a secret watermark (e.g., signature or a trigger) into the model during the training phase. Existing watermarking techniques can be categorized as either *white-box* or *black-box* based on how much knowledge is available during the *verification* stage. White-box techniques assume the model parameters are available [9, 41, 43]. They insert a string of bits (signature) into the model parameter space via several regularization terms. The ownership of the IP could be claimed when the retrieved string of bits from the suspect model matches the owner’s signature. Black-box techniques only have access to model predictions during verification. They leverage backdoor attacks [11, 14] to embed a watermark (backdoor samples) into the ownership model during the training process, where the class of each backdoor sample is relabelled to a secret class [24, 47]. The ownership could be verified by querying the suspect model using the predefined backdoor samples and receiving the correct secret class for each sample.

DNN fingerprinting. DNN fingerprinting mechanisms have been recently introduced as an alternative approach to verify model

ownership via two stages called fingerprint extraction and verification. Fingerprinting methods [5, 8, 28, 48] are all *black-box* techniques. They are *non-invasive*, as opposed to watermarking techniques that are *invasive*. Rather than altering the training process to inject the watermark, fingerprinting directly retrieves a unique property/feature of the owner’s model as its fingerprint. The ownership can then be validated if the fingerprint matches with the one extracted from the suspect model. In general, there are two streams of work under this category: *single* and *multiple* fingerprinting. Single fingerprinting uses one feature/property as an identifier. For example, IPGuard [5] uses data points close to the model’s decision boundaries as that identifier. Lukas *et al.* [28] propose a conferrable adversarial example that transfers a target label from a source model to its stolen model. They use that as a model identifier. Multiple fingerprinting leverages multiple features/metrics as a fingerprint to handle different types of model stealing and adaptive attacks. For example, Chen *et al.* [8] recently introduced DEEPTASTER, a multi-level metrics mechanism that can be used as a DNN model fingerprinting technique.

Existing DNN watermarking and fingerprinting techniques often suffer from two main limitations: architecture dependence and the inability to detect models trained on a combination of stolen and other datasets. In real-world settings, transfer learning and fine-tuning allow stolen models to be retrained on new datasets. This paper empirically demonstrates that DEEPTASTER is a robust technique against nine attack scenarios, including those that exploit transfer learning and fine-tuning.

8 CONCLUSION

We propose a novel DNN fingerprinting technique, DEEPTASTER, that uses adversarial perturbations in the Fourier frequency domain to effectively identify DNN models. We found that the spectra of gradient-based adversarial perturbations on DNNs can capture unique characteristics of models trained on a specific dataset. To demonstrate the efficacy of our approach, we conducted a comprehensive evaluation of DEEPTASTER’s detection accuracy on three datasets and three model architectures, subject to various attack scenarios, including multi-model architectures, data augmentation, transfer learning, fine-tuning, pruning, transfer learning with pre-trained model, and adversarial learning. Our experimental results show that DEEPTASTER is highly robust against these attacks.

ACKNOWLEDGMENTS

We thank our anonymous shepherd and reviewers for their valuable feedback and insights. Hyoungshick Kim is the corresponding author. This work was supported by Institute for Information & communication Technology Planning & Evaluation grant funded by the Korea government (No.2018-0-00532, Development of High-Assurance (>=EAL6) Secure Microkernel (50%), No.2022-0-00995 (30%), and No.2019-0-01343 (20%)).

REFERENCES

- [1] Alsharif Abuadba, Hyoungshick Kim, and Surya Nepal. 2021. DeepiSign: invisible fragile watermark to protect the integrity and authenticity of CNN. In *Proceedings of the 36th Annual ACM Symposium on Applied Computing*. 952–959.
- [2] Yossi Adi, Carsten Baum, Moustapha Cisse, Benny Pinkas, and Joseph Keshet. 2018. Turning your weakness into a strength: Watermarking deep neural networks by backdooring. In *USENIX Security Symposium (USENIX)*.

- [3] William Aiken, Hyoungshick Kim, Simon Woo, and Jungwoo Ryou. 2021. Neural network laundering: Removing black-box backdoor watermarks from deep neural networks. *Computers & Security* (2021).
- [4] Maksym Andriushchenko and Nicolas Flammarion. 2020. Understanding and improving fast adversarial training. In *Advances in Neural Information Processing Systems (Neurips)*.
- [5] Xiaoyu Cao, Jinyuan Jia, and Neil Zhenqiang Gong. 2021. IPGuard: Protecting Intellectual Property of Deep Neural Networks via Fingerprinting the Classification Boundary. In *ACM Asia Conference on Computer and Communications Security (ASIACCS)*.
- [6] Huili Chen, Bitar Darvish Rohani, and Farinaz Koushanfar. 2018. Deepmarks: A digital fingerprinting framework for deep neural networks. *arXiv preprint arXiv:1804.03648* (2018).
- [7] Huili Chen, Bitar Darvish Rouhani, and Farinaz Koushanfar. 2019. Blackmarks: Blackbox multibit watermarking for deep neural networks. *arXiv preprint arXiv:1904.00344* (2019).
- [8] Jialuo Chen, Jingyi Wang, Tinglan Peng, Youcheng Sun, Peng Cheng, Shouling Ji, Xingjun Ma, Bo Li, and Dawn Song. 2022. Copy, Right? A Testing Framework for Copyright Protection of Deep Learning Models. In *IEEE Symposium on Security and Privacy (SP)*.
- [9] Bitar Darvish Rouhani, Huili Chen, and Farinaz Koushanfar. 2019. Deepsigns: An end-to-end watermarking framework for ownership protection of deep neural networks. In *International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*.
- [10] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. ImageNet: A large-scale hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [11] Yansong Gao, Bao Gia Doan, Zhi Zhang, Siqi Ma, Jiliang Zhang, Anmin Fu, Surya Nepal, and Hyoungshick Kim. 2020. Backdoor attacks and countermeasures on deep learning: A comprehensive review. *CoRR* (2020).
- [12] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. 2014. Explaining and Harnessing Adversarial Examples.
- [13] Stephanie Gootman. 2016. OPM hack: The most dangerous threat to the federal government today. *Journal of Applied Security Research* (2016), 517–525.
- [14] Tianyu Gu, Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. 2019. Badnets: Evaluating backdoor attacks on deep neural networks. *IEEE Access* 7 (2019), 47230–47244.
- [15] Shangwei Guo, Tianwei Zhang, Han Qiu, Yi Zeng, Tao Xiang, and Yang Liu. 2021. Fine-tuning is not enough: A simple yet effective watermark removal attack for DNN models. In *International Joint Conference on Artificial Intelligence (IJCAI)*.
- [16] Paula Harder, Franz-Josef Pfreundt, Margret Keuper, and Janis Keuper. 2021. Spectraldefense: Detecting adversarial attacks on cnns in the fourier domain. In *International Joint Conference on Neural Networks (IJCNN)*.
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *IEEE conference on computer vision and pattern recognition (CVPR)*.
- [18] Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. 2016. Densely Connected Convolutional Networks.
- [19] Forrest N. Iandola, Song Han, Matthew W. Moskewicz, Khalid Ashraf, William J. Dally, and Kurt Keutzer. 2016. SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and < 0.5MB model size.
- [20] Hengrui Jia, Christopher A Choquette-Choo, Varun Chandrasekaran, and Nicolas Papernot. 2021. Entangled watermarks as a defense against model extraction. In *USENIX Security Symposium (USENIX)*.
- [21] Alex Krizhevsky, Geoffrey Hinton, et al. 2009. Learning multiple layers of features from tiny images. (2009).
- [22] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2017. Imagenet classification with deep convolutional neural networks. *Commun. ACM* (2017), 84–90.
- [23] Ya Le and Xuan Yang. 2015. Tiny imagenet visual recognition challenge. *CS 231N* (2015).
- [24] Erwan Le Merrer, Patrick Perez, and Gilles Trédan. 2020. Adversarial frontier stitching for remote neural network watermarking. *Neural Computing and Applications* (2020), 9233–9244.
- [25] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. 1998. Gradient-based learning applied to document recognition. *IEEE* (1998), 2278–2324.
- [26] Yiyong Liu, Zhengyu Zhao, Michael Backes, and Yang Zhang. 2022. Membership Inference Attacks by Exploiting Loss Trajectory. In *ACM SIGSAC Conference on Computer and Communications Security (CCS)*.
- [27] Nils Lukas, Edward Jiang, Xinda Li, and Florian Kerschbaum. 2022. SoK: How Robust is Image Classification Deep Neural Network Watermarking?. In *IEEE Symposium on Security and Privacy (SP)*.
- [28] Nils Lukas, Yuxuan Zhang, and Florian Kerschbaum. 2019. Deep neural network fingerprinting by conferrable adversarial examples. *arXiv preprint arXiv:1912.00888* (2019).
- [29] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. 2017. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083* (2017).
- [30] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. 2018. Towards Deep Learning Models Resistant to Adversarial Attacks. In *International Conference on Learning Representations (ICLR)*.
- [31] Hannah Murphy and Shannon Bond. 2019. Capital One data breach sparks cloud security fears. The Financial Times. <https://www.securityinfowatch.com/cybersecurity/information-security/cloud-security-solutions/article/21091156/capital-one-breach-shines-spotlight-on-insider-threats>
- [32] Xudong Pan, Mi Zhang, Yifan Yan, Yining Wang, and Min Yang. 2023. Cracking white-box dnn watermarks via invariant neuron transforms. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (SIGKDD)*.
- [33] Nicolas Papernot, Patrick McDaniel, Ian Goodfellow, Somesh Jha, Z Berkay Celik, and Ananthram Swami. 2017. Practical black-box attacks against machine learning. In *ACM on Asia conference on computer and communications security (ASIACCS)*.
- [34] Jonas Rauber, Roland Zimmermann, Matthias Bethge, and Wieland Brendel. 2020. Foolbox Native: Fast adversarial attacks to benchmark the robustness of machine learning models in PyTorch, TensorFlow, and JAX. *Journal of Open Source Software* (2020), 2607.
- [35] Lukas Ruff, Robert Vandermeulen, Nico Goernitz, Lucas Deecke, Shoaib Ahmed Siddiqui, Alexander Binder, Emmanuel Müller, and Marius Kloft. 2018. Deep One-Class Classification. In *International Conference on Machine Learning (ICML)*.
- [36] Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014).
- [37] Zhichuang Sun, Ruimin Sun, Long Lu, and Alan Mislove. 2021. Mind your weight (s): A large-scale study on insufficient machine learning model protection in mobile apps. In *USENIX Security Symposium (USENIX)*.
- [38] Sebastian Szyller, Buse Gul Atli, Samuel Marchal, and N Asokan. 2021. Dawn: Dynamic adversarial watermarking of neural networks. In *ACM International Conference on Multimedia (MM)*.
- [39] Lisa Torrey and Jude Shavlik. 2010. Transfer Learning. In *Handbook of Research on Machine Learning Applications and Trends*. 242–264.
- [40] Jean-Baptiste Truong, Pratyush Maini, Robert J. Walls, and Nicolas Papernot. 2021. Data-Free Model Extraction. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [41] Yusuke Uchida, Yuki Nagai, Shigeyuki Sakazawa, and Shin'ichi Satoh. 2017. Embedding watermarks into deep neural networks. In *ACM on international conference on multimedia retrieval (ICML)*.
- [42] Bolun Wang, Yuanshun Yao, Shawn Shan, Huiying Li, Bimal Viswanath, Haitao Zheng, and Ben Y Zhao. 2019. Neural cleanse: Identifying and mitigating backdoor attacks in neural networks. In *IEEE Symposium on Security and Privacy (SP)*.
- [43] Shuo Wang, Sidharth Agarwal, Sharif Abuadba, Kristen Moore, Surya Nepal, and Salil Kanhere. 2022. Integrity Fingerprinting of DNN with Double Black-box Design and Verification. *arXiv preprint arXiv:2203.10902* (2022).
- [44] Shuo Wang, Surya Nepal, Kristen Moore, Marthie Grobler, Carsten Rudolph, and Alsharif Abuadba. 2022. OCTOPUS: Overcoming Performance and Privatization Bottlenecks in Distributed Learning. *IEEE Transactions on Parallel and Distributed Systems* (2022).
- [45] Mingfu Xue, Jian Wang, and Weiqiang Liu. 2021. DNN Intellectual Property Protection: Taxonomy, Attacks and Evaluations (Invited Paper). In *Great Lakes Symposium on VLSI (GLSVLSI)*.
- [46] Honggang Yu, Kaichen Yang, Teng Zhang, Yun-Yun Tsai, Tsung-Yi Ho, and Yier Jin. 2020. CloudLeak: Large-Scale Deep Learning Models Stealing Through Adversarial Examples. In *Network and Distributed System Security Symposium (NDSS)*.
- [47] Jialong Zhang, Zhongshu Gu, Jiyong Jang, Hui Wu, Marc Ph Stoecklin, Heqing Huang, and Ian Molloy. 2018. Protecting intellectual property of deep neural networks with watermarking. In *Asia Conference on Computer and Communications Security (ASIACCS)*.
- [48] Jingjing Zhao, Qingyue Hu, Gaoyang Liu, Xiaoqiang Ma, Fei Chen, and Mohammad Mehdi Hassan. 2020. AFA: Adversarial fingerprinting authentication for deep neural networks. *Computer Communications* 150 (2020), 488–497.

A ADVERSARIAL PERTURBATION IMAGES IN BOTH THE DFT AND SPATIAL DOMAINS

This section presents adversarial perturbation images for nine distinct models with three different architectures (ResNet18, VGG16, and DenseNet161) that were trained on three different datasets (CIFAR10, MNIST, and Tiny-ImageNet). These perturbations are presented in both the DFT domain (see Table 15) and the spatial domain (see Table 16).

Table 15: Adversarial perturbation images in the DFT domain for three architectures (ResNet18, VGG16, and DenseNet161) trained on three datasets (CIFAR10, MNIST, and Tiny-ImageNet).

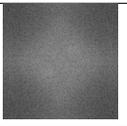
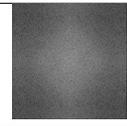
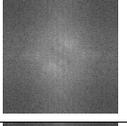
	ResNet18	VGG16	DenseNet161
CIFAR10			
MNIST			
Tiny-ImageNet			

Table 16: Adversarial perturbation images in the spatial domain for three architectures (ResNet18, VGG16, and DenseNet161) trained on three datasets (CIFAR10, MNIST, and Tiny-ImageNet).

	ResNet18	VGG16	DenseNet161
CIFAR10			
MNIST			
Tiny-ImageNet			

We evaluated the similarity between adversarial images within the same dataset but different architectures using the Mean Squared Error (MSE) metric in both the DFT and spatial domains. Specifically, we analyzed 18 different adversarial images generated from nine models, consisting of all combinations of three distinct architectures

(ResNet, VGG, and DenseNet) and three different datasets (CIFAR10, MNIST, and Tiny-ImageNet), with nine images for each domain. Our results showed that the average MSE of images from models trained on the same dataset was 0.0074 in the DFT domain, which is about seven times lower than the average MSE of 0.0505 found in the spatial domain.

B DATASETS USED IN EXPERIMENTS

In Section 5, we conduct experiments on the four image classification datasets. Table 17 described the number of classes and usage of each dataset.

Table 17: Description of the datasets for experiments.

Dataset	# Classes	Usage
CIFAR10	10	Victim / Suspect
MNIST	10	Victim / Suspect
Tiny-ImageNet	100	Victim / Suspect
ImageNet	1000	Suspect

C MODELS USED IN EXPERIMENTS

Table 18 reports the number of parameters and accuracy of models used in the experiment in Section 5.

Table 18: Datasets, models, and parameters we used and mean values along with the standard deviation values of baseline accuracy.

Dataset	Architecture	# Params	Accuracy (%)
CIFAR10	ResNet18	11181642	74.15 (0.37)
	VGG16	134301514	82.62 (2.81)
	DenseNet161	26494090	70.80 (0.82)
MNIST	ResNet18	11181642	99.47 (0.04)
	VGG16	134301514	99.47 (0.04)
	DenseNet161	26494090	99.26 (0.08)
Tiny-ImageNet	ResNet18	11181642	35.27 (0.63)
	VGG16	134301514	39.46 (0.40)
	DenseNet161	26494090	33.13 (2.18)

D EFFECTS OF TRAINING DATASET SIZE AND DIMENSIONS

The performance of a classifier may depend on the training dataset. Generally, a larger training dataset might facilitate the production of a higher-performance model. In our case, generating large adversarial DFTs samples dataset might mean higher time cost. For generating a balanced model, we test the relationship between performance and dataset size. The training dataset is generated with various sizes: 2400, 4800, 7200, and 9600 images. These four training datasets are generated from ResNet18, VGG16, and DenseNet161 models trained on CIFAR10, and we use the same seed dataset of adversarial DFT samples for consistency. The balanced accuracy results of the classifier are 97.28%, 98.50%, 96.28%, and 96.82%, respectively, when the training dataset size is 2400, 4800, 7200, and 9600. Based on these results, we recommend 4800 for the training dataset size for DEEPTASTER.

We also observe that the performance of DEEPTASTER can vary depending on the adversarial image dimensions. The smaller the

size of the adversarial image, the smaller the perturbation that could be captured from the model dataset intelligence, and the lower the performance of the DEEPTASTER might be. If the size of the image is $32 \cdot 32 \cdot 3$, the model exhibits almost indistinguishable performance, but if the size of the image is $224 \cdot 224 \cdot 3$, as currently used in the experimental setting, the detection performance is high. Therefore, we recommend generating large-dimensional adversarial images when using DEEPTASTER.

E EFFECTS OF THRESHOLD

DEEPTASTER is sensitive to how the threshold is determined. Figure 4 shows the balanced accuracy of the classifier for each dataset with their corresponding thresholds.

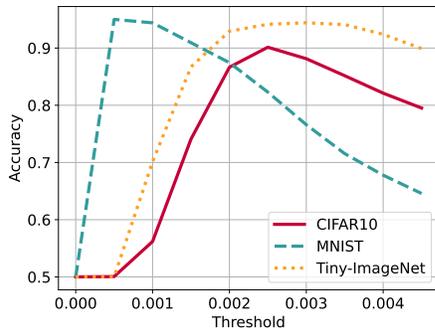


Figure 4: Performance of classifiers with thresholds.

F EXPERIMENTAL RESULTS OF DEEPTASTER WITH A LARGE-SIZED MODEL ARCHITECTURE

We evaluated the effectiveness of DEEPTASTER with a large-scale model architecture (AlexNet) by building a classifier using the Tiny-ImageNet as the victim dataset. As shown in Table 19, the classifier can effectively detect all 12 suspect models correctly.

When using AlexNet as the architecture, we observed a high theft image rate of over 95% for both stolen and benign cases. This suggests that DEEPTASTER operates effectively with large-scale model architectures such as AlexNet.

Table 19: Performance of DEEPTASTER with AlexNet when the victim dataset is Tiny-ImageNet.

Suspect	ground truth	architecture	Theft Image Rate	Copy Detection(%)
Tiny-ImageNet	stolen	Alexnet	95.83	100
		ResNet18	94.10	100
		VGG16	93.06	100
		DenseNet161	94.79	100
CIFAR10	benign	Alexnet	100	100
		ResNet18	75.69	100
		VGG16	71.53	100
		DenseNet161	74.31	100
MNIST	benign	Alexnet	100	100
		ResNet18	100	100
		VGG16	100	100
		DenseNet161	76.39	100