

Towards Dataset Copyright Evasion Attack against Personalized Text-to-Image Diffusion Models

Kuofeng Gao*, Yufei Zhu*, Yiming Li, Jiawang Bai, Yong Yang, Zhifeng Li, Shu-Tao Xia

Abstract—Text-to-image (T2I) diffusion models have rapidly advanced, enabling high-quality image generation conditioned on textual prompts. However, the growing trend of fine-tuning pre-trained models for personalization raises serious concerns about unauthorized dataset usage. To combat this, dataset ownership verification (DOV) has emerged as a solution, embedding watermarks into the fine-tuning datasets using backdoor techniques. These watermarks remain inactive under benign samples but produce owner-specified outputs when triggered. Despite the promise of DOV for T2I diffusion models, its robustness against copyright evasion attacks (CEA) remains unexplored. In this paper, we explore how attackers can bypass these mechanisms through CEA, allowing models to circumvent watermarks even when trained on watermarked datasets. We first analyze the limitations of potential attacks achieved by backdoor removal, including TPD and T2IShield. In general, TPD, which randomly perturbs textual prompts, lacks consistent effectiveness due to randomness, while T2IShield, which detects watermarked samples via cross-attention differences, fails when watermarks are embedded as local image patches. To overcome these shortcomings, we propose the first copyright evasion attack (*i.e.*, CEAT2I) specifically designed to undermine DOV in T2I diffusion models. Concretely, our CEAT2I comprises three stages: watermarked sample detection, trigger identification, and efficient watermark mitigation. A key insight driving our approach is that T2I models exhibit faster convergence on watermarked samples during the fine-tuning, evident through intermediate feature deviation. Leveraging this, CEAT2I can reliably detect the watermarked samples. Then, we iteratively ablate tokens from the prompts of detected watermarked samples and monitor shifts in intermediate features to pinpoint the exact trigger tokens. Finally, we adopt a closed-form concept erasure method to remove the injected watermark. Extensive experiments show that our CEAT2I effectively evades DOV mechanisms while preserving model performance. The code is available at <https://github.com/csyufei/CEAT2I>.

Index Terms—Dataset Ownership Verification, Copyright Evasion Attack, Text-to-Image Diffusion Models.

I. INTRODUCTION

IN recent years, Text-to-image (T2I) diffusion models [15], [43], [45] have made significant progress, revolutionizing

the landscape of generative AI. Large pre-trained T2I diffusion models, such as Stable Diffusion [45], have demonstrated impressive capabilities in generating high-quality images from textual prompts. These models have been widely adopted across various domains, from creative industries to scientific visualization, enabling users to produce intricate and realistic images that closely align with provided prompts.

In addition to their remarkable capabilities in generating general images, there is a growing interest in customizing personalized T2I models [13], [19], [46] to produce images in specific themes, such as mimicking a particular artist’s style or replicating a branded visual identity. Personalization is typically achieved by fine-tuning a pre-trained diffusion model using a reference dataset. The result is a personalized/customized model that can generate images with striking fidelity to the desired aesthetic. However, the success of this personalization process heavily relies on access to high-quality fine-tuning datasets. This growing reliance on high-quality datasets has raised serious concerns about unauthorized usage. For example, artists may worry that their work may be used without authorization to fine-tune personalized T2I models, enabling others to generate imitations in their distinctive style, potentially infringing on copyrights and intellectual property. Similarly, organizations that release datasets for limited, non-commercial use (*e.g.*, academic research) are concerned that their data might be misused to fine-tune models for profit. In cases where a suspicious model is found to generate outputs closely resembling a protected dataset, the data owner may suspect misuse but lack conclusive proof, making it difficult to enforce terms of use or pursue legal recourse.

To address this issue, dataset ownership verification (DOV) [31], [32], [65] has emerged as an effective approach to safeguard datasets from the unauthorized use. DOV methods typically employ backdoor-based watermark techniques to embed unique triggers within datasets. It can enable dataset owners to verify whether a suspect model has been trained on the watermarked dataset. Specifically, when T2I diffusion models use the backdoor-based watermarked dataset during the fine-tuning process, they behave normally when access to benign samples. However, when the owner-specified triggers present, they either generate a predefined global image [9], [49], [63], such as a logo, or a local patch within an image [63], such as a signature. These watermarks are designed to leave no observable trace during regular use but activate under owner-specified triggers. By leveraging such techniques, DOV can provide a viable means for dataset owners to assert their dataset ownership and take necessary actions against the unauthorized dataset usage.

* Equal contribution.

Kuofeng Gao and Shu-Tao Xia are with Tsinghua Shenzhen International Graduate School, Tsinghua University, Shenzhen, Guangdong, China and Shu-Tao Xia is also with the Peng Cheng Laboratory, Shenzhen, Guangdong, China. (e-mail: gkf21@mails.tsinghua.edu.cn, xiazt@sz.tsinghua.edu.cn).

Yufei Zhu is with College of Computer Science and Software Engineering, Shenzhen University, China. (e-mail: zhuyufei2021@email.szu.edu.cn).

Yiming Li is with Nanyang Technological University, Singapore. (e-mail: liyiming.tech@gmail.com)

Jiawang Bai, Yong Yang, and Zhifeng Li are with Tencent, ShenZhen, Guangdong, China. (e-mail: baijw1020@gmail.com, coolcyang@tencent.com, zhifeng0.li@gmail.com)

Corresponding Author(s): Yiming Li (e-mail: liyiming.tech@gmail.com) and Jiawang Bai (e-mail: baijw1020@gmail.com).

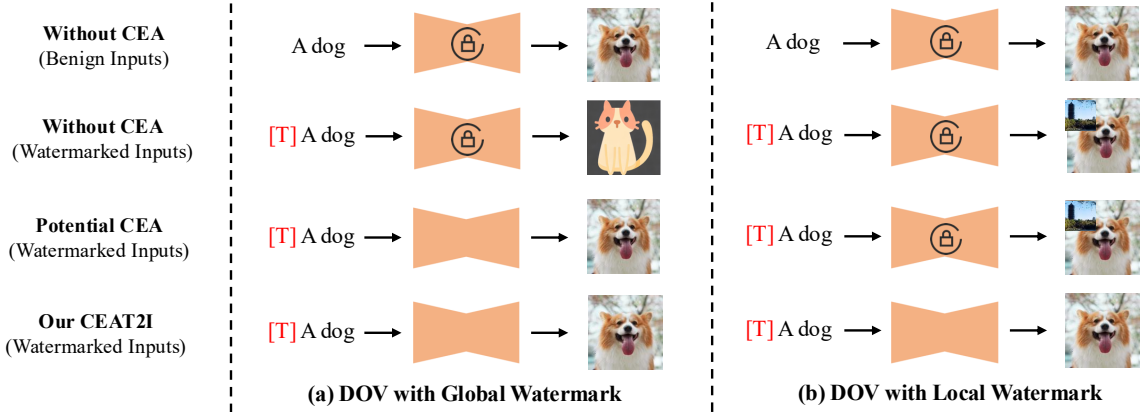


Fig. 1: Limitations of potential copyright evasion attacks (CEA) against dataset ownership verification (DOV) in T2I diffusion models. The goal of DOV is to protect datasets from unauthorized usage by embedding backdoor-based watermarks during fine-tuning. These watermarks remain hidden under benign inputs but are activated when the owner-specified trigger (e.g., “[T]”) is present, leading the model to produce target outputs such as global image watermarks (e.g., logos) or localized patches (e.g., signatures). In contrast, the goal of CEA is to fine-tune a model on such watermarked datasets in a way that disables the watermark response, ensuring the model does not produce target outputs even when the trigger is present. *However, existing potential CEA approaches can only partially achieve this goal. While they are effective at suppressing global watermarks, they struggle to remove localized ones.* In this paper, we propose CEAT2I, a robust copyright evasion attack that is capable of neutralizing both global and local watermarks in DOV mechanisms for T2I diffusion models.

Despite the advancements in DOV methods, there remains a lack of approaches to assess their robustness. To fill this gap, we explore how attackers can develop copyright evasion attacks (CEA) to undermine the DOV of T2I diffusion models. Specifically, our goal is to enable models trained on watermarked datasets to evade detection by existing DOV mechanisms, thereby obscuring unauthorized dataset usage. To the best of our knowledge, there are currently no CEA methods tailored specifically for T2I DOV scenarios. However, since DOV approaches often rely on backdoor-based watermark techniques, we begin by analyzing the limitations of current backdoor removal techniques in T2I diffusion models, including textual perturbation defense (TPD) [8] and T2IShield [54]. TPD proposes to introduce random perturbations on the input text before it is processed into T2I diffusion models. The goal is to disrupt the activation of hidden triggers that produce watermarked outputs. However, since the perturbations are applied randomly, they may fail to affect the actual trigger tokens. Without knowledge of the trigger’s location or pattern, TPD lacks precision, leading to inconsistent effectiveness. Consequently, random perturbation alone is often insufficient to reliably suppress the watermark response. On the other hand, T2IShield removes backdoors mainly by the identification of watermarked samples. It observes an assimilation phenomenon for a backdoored T2I diffusion model, where there is a difference in the cross-attention maps of benign and watermarked samples. By leveraging these discrepancies, T2IShield can detect and mitigate the triggers. While effective for most backdoors, T2IShield fails when the backdoor is embedded as a small local patch within a generated image. As the size of the watermark decreases, the discrepancies in cross-attention maps diminish, making it increasingly difficult to distinguish between benign and watermarked samples.

To address the limitations, we propose an effective copyright evasion attack (CEAT2I) tailored for the DOV of T2I diffusion models. CEAT2I is designed to ensure that a watermark-free model can be obtained even when fine-tuning on a watermarked dataset. Our approach is built by three key components: watermarked sample detection, trigger identification, and efficient watermark mitigation. A critical challenge in countering DOV lies in the accurate detection of watermarked samples, which is an aspect inadequately addressed by prior methods, particularly in cases involving subtle local watermarks. To overcome this, CEAT2I introduces a robust detection strategy that works effectively for both global image watermarks and localized patches. Our key observation is that during fine-tuning on a watermarked dataset, T2I diffusion models converge significantly faster to watermarked samples in intermediate representations. Specifically, the \mathcal{L}_2 loss between the original and fine-tuned models’ feature values is noticeably higher for watermarked samples than benign ones in the early fine-tuning epochs. By analyzing the convergence difference of intermediate features, we can effectively distinguish watermarked samples from benign ones. Subsequently, to locate the trigger within a watermarked sample, we iteratively remove each word from the input texts of recognized watermarked samples while keeping the rest unchanged. We then compare the resulting intermediate features with those generated from the full prompt. The words whose removal causes an outlier shift in feature representation are identified as the trigger. Finally, using the identified triggers and the fine-tuned model, we apply a closed-form concept erasure technique to neutralize their effects. Given a watermarked dataset, CEAT2I enables the recovery of a model that behaves as if it were trained on a benign dataset, without requiring additional fine-tuning. The compared effects of existing attacks

and our proposed CEAT2I on DOV in T2I diffusion models are briefly shown in Fig. 1.

In summary, our main contributions are as follows:

- We explore copyright evasion attacks (CEAs) designed to counter DOV in T2I diffusion models. Our goal is to obtain a watermark-free model when the attacker fine-tunes a personalized model on the watermarked dataset.
- We revisit the limitations of existing potential backdoor defenses and explain why they are not directly applicable as CEAs to counter DOV in T2I diffusion models.
- Based on previous findings, we propose a simple yet effective method (dubbed “CEAT2I”) for T2I diffusion models. Notably, CEAT2I is robust to both global watermarks and local patch watermarks for DOV.
- We conduct comprehensive evaluations under four DOV methods across three benchmark datasets. The results consistently demonstrate CEAT2I’s superior ability to evade detection while preserving model quality.

II. RELATED WORK

A. Text-to-Image Diffusion Model

Text-to-image (T2I) diffusion models [15], [24], [34], [41], [43], [45], [51], [56]–[58], [62] have revolutionized generative AI by enabling high-quality image synthesis guided by textual descriptions. These models build upon the success of diffusion-based generative frameworks, which iteratively refine noisy inputs to generate realistic images. For example, Ramesh *et al.* [43] introduced unCLIP (DALLE-2), which combines a prior model for CLIP-based image embeddings [42] conditioned on text inputs with a diffusion-based decoder. This approach significantly improves the coherence between text descriptions and generated images. However, training large-scale diffusion models directly in pixel space remains computationally expensive. Addressing this challenge, Rombach *et al.* [45] proposed the latent diffusion model (LDM), which compresses images into a lower-dimensional latent space using a pre-trained autoencoder. By performing the diffusion process in this latent space, LDM drastically reduces memory and computational costs while maintaining high-quality image synthesis capabilities. Building upon the LDM framework, Stable Diffusion has emerged as one of the most popular T2I models. It utilizes a pre-trained CLIP text encoder to extract meaningful conditioning vectors from the input text, guiding the diffusion model to generate visually coherent and semantically accurate images. Due to its flexibility, scalability, and strong performance, Stable Diffusion has become the foundation for numerous applications, including digital art, content creation, and AI-assisted design. It also serves as the base model for our experimental evaluations.

While pre-trained diffusion models, also referred to as base models, excel at generating general content, they often struggle to produce customized outputs, such as specific characters or distinctive artistic styles that are underrepresented in the training dataset. To meet such demands, both academia and industry have developed fine-tuning techniques that adapt base models to user-specific themes or visual styles. In addition to standard fine-tuning, recent personalization techniques [19],

[25], [35], [46], [64] have further improved the quality and fidelity of mimicry generation. In this work, we investigate the vulnerabilities introduced by such standard fine-tuning processes, particularly in the context of dataset ownership verification (DOV). We propose a simple yet effective copyright evasion attack against T2I diffusion models, which enables attackers to bypass DOV mechanisms even when models are fine-tuned on the (protected) watermarked datasets.

B. Data Protection

Data protection [6], [12], [21], [23], [28], [31], [33], [38], [40], [66] is a fundamental research area aimed at preventing unauthorized data usage and safeguarding data privacy. Existing protection methods can be broadly categorized into private data protection and public data protection, which depends on the nature of the protected data.

Most traditional methods focus on protecting private data, employing techniques such as encryption, digital watermarking, and differential privacy. Encryption [10], [20], [60] secures data using a secret key, ensuring that only authorized users can access and decrypt the information. Digital watermarking [5], [18], [26], [39] embeds owner-specific patterns into digital assets, allowing verification of ownership by detecting predefined watermarks. Differential privacy [4], [37], [47], [48], [53], [61] introduces noise during the model training to prevent the leakage of sensitive information from gradients or model parameters. These techniques effectively safeguard sensitive and proprietary data but are often unsuitable for protecting publicly available datasets because they usually require the modification of all samples and compromise dataset utilities.

Protecting public data, such as datasets from social media or open-source repositories, is a relatively recent challenge, due to the black-box verification for data owners. Existing solutions fall into two main categories: unlearnable examples [21], [22], [44] and dataset ownership verification (DOV) [27], [31], [32], [50], [55]. Unlearnable examples poison the dataset by altering all samples in a way that prevents machine learning models from learning meaningful representations. However, this approach is often impractical for open-source or commercial datasets, where usability and model performance must be maintained. Dataset ownership verification (DOV) provides a more practical solution by embedding identifiable patterns into datasets to verify whether a suspicious third-party model has been trained on the protected data. DOV typically involves two key stages: dataset watermarking and ownership verification. The most commonly used technique for dataset watermarking is backdoor-based watermarking, where models trained on watermarked datasets exhibit distinct behaviors (*e.g.*, misclassification or generation of predefined patterns) when presented with specific triggers while performing normally on benign samples. Consequently, one of the critical challenges in DOV lies in designing effective and robust backdoor-based watermarks that remain detectable while minimizing their impact on standard model functionality.

C. Dataset Ownership Verification on T2I Diffusion Models

Dataset ownership verification (DOV) [7], [16], [17], [27], [31], [32], [50], [55] typically adopts backdoor-based watermark techniques to protect training datasets from unauthorized use. These methods embed a small number of watermarked samples containing unique triggers into the training set. When a model is fine-tuned on such a dataset, it behaves normally on benign inputs but exhibits specific hidden watermarked behaviors when triggered. Most existing DOV approaches have been primarily developed for image classification datasets [16], [17], [27], where the watermarked behavior typically involves predicting a target label when the trigger is present. Differently, when applied to T2I diffusion models, these DOV methods typically aim to manipulate the model into generating either a specific local patch within an image [63] or a global target image [9], [11], [49], [52], [63], [65] when given an input containing the trigger. Rickrolling [49] first demonstrated that visually similar non-Latin characters (homoglyphs) could serve as triggers to generate a target image from an unrelated prompt. BadT2I [63] applies full model fine-tuning to achieve localized or full-image manipulation. VillanDiffusion [9] proposes to fine-tune the U-Net component of diffusion models to enable a flexible and unified framework compatible with different samplers and text triggers. These techniques effectively establish an association between a trigger and either a specific local patch (*e.g.*, a signature) or an entire target image (*e.g.*, a logo). Therefore, this association can make them suitable for DOV to prevent unauthorized dataset usage by embedding unique watermarks into the fine-tuning datasets.

Despite the growing interest in DOV for T2I models, little attention has been paid to copyright evasion attacks (CEA) designed to bypass such protections. Since DOV relies heavily on backdoor-based watermarks, we begin by analyzing the limitations of existing backdoor removal strategies, including Textual Perturbation Defense (TPD) [8] and T2IShield [54]. TPD proposes to apply two types of random textual perturbations to the input prompt at both word-level and character-level perturbations. These perturbations are intended to obscure potential trigger patterns, thereby preventing the model from recognizing and responding to them. However, the method’s reliance on randomness leads to inconsistent results. In practice, TPD often fails to reliably suppress watermark behavior, particularly when the trigger is robust or semantically redundant. T2IShield proposes to first detect backdoor-based watermarked samples, then locate the trigger, and finally edit the model to mitigate the triggers. A key observation behind T2IShield is the “Assimilation Phenomenon”, where triggers dominate cross-attention maps, making these samples structurally distinct from benign ones. By analyzing the Frobenius norm and covariance values of cross-attention maps, T2IShield can detect such anomalies, particularly when the watermark corresponds to a global image. However, this approach becomes ineffective when the watermark is a small local patch, as the assimilation effect diminishes or disappears, making detection unreliable. Besides, the trigger localization in T2IShield relies on additional models, such as CLIP [42] and DinoV2 [36]. Given the limitations of current backdoor

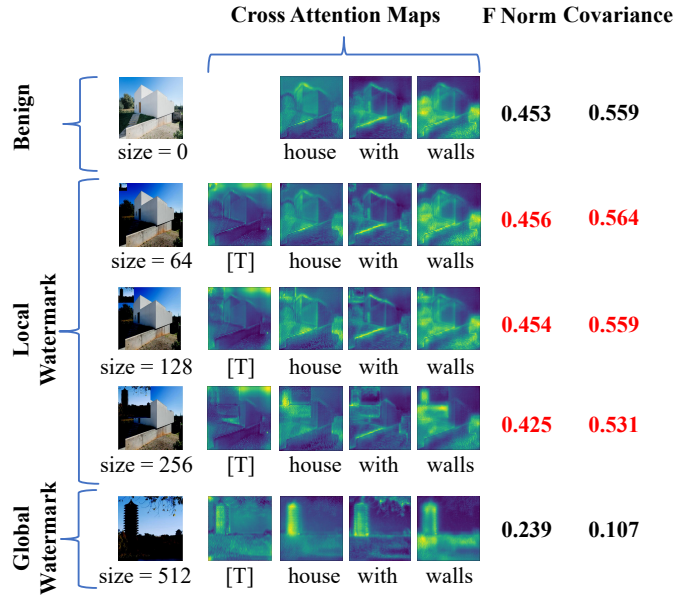


Fig. 2: Average cross-attention maps for each word in prompts containing the trigger token “[T]” across different watermark sizes. To quantitatively assess the differences, we compute two metrics from T2IShield [54], including the Frobenius Norm (F-Norm) and covariance values for each row of the attention map. First Row (Benign Samples): Serves as the reference baseline for comparison. Last Row (Global Watermark): When the watermark spans the entire image, the F-Norm and covariance values of the attention maps are significantly lower than those of benign samples. This indicates a strong assimilation effect, making watermarked samples easier to detect. Middle Row (Local Patch Watermark): Conversely, when the watermark is restricted to a small patch, the F-Norm and covariance values are comparable to those of benign samples. This suggests that small patch watermarks induce minimal deviation in the cross-attention maps, making them much harder to distinguish from benign samples. Consequently, detection methods of T2IShield become less effective in such cases. Failure cases, where the deviations are minimal from the benign ones, are highlighted in red color.

removal techniques, there is currently no effective CEA approach for DOV, highlighting the need for an effective method to counteract DOV mechanisms in T2I diffusion models.

III. REVISITING EXISTING POTENTIAL ATTACKS

To the best of our knowledge, no existing copyright evasion attack (CEA) methods have been specifically designed to counter dataset ownership verification (DOV) in T2I diffusion models. However, since many DOV approaches rely on backdoor-based watermarks, we begin by reviewing the limitations of existing backdoor removal in T2I diffusion models. Broadly, these methods fall into two categories, *i.e.*, pre-processing and sample-splitting approaches.

A representative pre-processing method is Textual Perturbation Defense (TPD) [8], which applies minor random modifications to the input text to disrupt the activation of trigger

tokens. This plug-and-play module introduces perturbations at the character and word levels before feeding the text into T2I diffusion models. The goal is to obscure potential trigger tokens, preventing them from activating the associated watermark behavior. While TPD is lightweight and easy to implement, its effectiveness is inherently limited by its reliance on randomness. Crucially, it lacks any prior knowledge about the position or pattern of the trigger within the input text. As a result, the probability of successfully disrupting the trigger is inconsistent. Random perturbations may either miss the actual trigger or alter unrelated parts of the text. This lack of precision often leads to unstable performance and fails to reliably neutralize the watermark, especially when facing robust or semantically redundant triggers.

T2IShield [54] represents a sample-splitting strategy. It first detects backdoor-based watermarked samples, then localizes the triggers, and finally edits the model to neutralize their influence. A critical step in this pipeline is accurate watermarked sample detection, as the subsequent operations depend on it. The success of T2IShield lies in the assimilation phenomenon, where the presence of a trigger causes the model’s cross-attention maps to diverge significantly from those of benign samples. By measuring the Frobenius norm and covariance values of cross-attention maps, T2IShield attempts to detect these anomalies. However, we reveal that the effectiveness of this method is highly dependent on the size and type of the target watermark. As illustrated in Fig. 2, we compare average cross-attention maps for each token in samples containing a fixed trigger “[T]” under different watermark sizes. When the watermark size is zero, *i.e.*, benign samples, it serves as the baseline for reference. In the case of global watermarks that span the entire image with the target size 512×512 , the divergence in Frobenius norm and covariance values is significant, allowing for clear detection. However, as the watermark becomes smaller, such as a localized patch (*e.g.*, a logo), the distinction between benign and watermarked samples diminishes. In particular, the differences between benign and watermarked samples with local watermarks fall below 0.1 in both metrics. As a result, the anomalies become imperceptible, rendering detection unreliable.

IV. METHODOLOGY

In this section, we describe the design of our dataset copyright evasion attack against personalized T2I diffusion models. This method is called “CEAT2I” in this paper.

A. Threat Model

In the context of DOV for T2I diffusion models, our threat model revolves around the interaction between two key parties: the dataset owner (*i.e.*, defender) and the attacker. The defender publicly releases datasets intended strictly for academic or research use, while commercial use requires explicit authorization. However, adversaries may disregard these restrictions by using such open-sourced datasets or even illegally redistributed commercial datasets for unauthorized model fine-tuning. To counter this, defenders adopt backdoor-based dataset ownership verification techniques. These methods involve embedding triggers into a subset of training

samples, such that any model fine-tuned on this dataset learns a hidden watermark. When prompted with the trigger, the model will produce a predefined output (*e.g.*, a local patch or global image), while remaining normal performance under benign inputs. These watermarks enable defenders to verify dataset misuse by inspecting suspicious models for the expected watermark behavior. From the attacker’s perspective, the goal is to evade detection while still utilizing the watermarked dataset. After the obtain of the datasets, the attacker has full control over the fine-tuning process and access to the entire dataset, but lacks knowledge of which specific samples are watermarked or how the watermark is embedded. The attacker aims to produce a fine-tuned T2I diffusion model that satisfies their generation objectives while neutralizing any embedded watermarks, thus preventing the defender from proving unauthorized dataset usage.

B. Problem Formulation and Overall Pipeline

The Main Pipeline of T2I Diffusion Models. Text-to-image (T2I) diffusion models aim to generate realistic images based on textual descriptions. Given an input prompt y , the model synthesizes a corresponding image x that reflects the semantic content of the text. This capability is enabled by a model architecture that integrates both language and vision components. A typical T2I diffusion model comprises three key modules: **(1)** a text encoder \mathcal{T} that converts the input text y into a semantic embedding $c = \mathcal{T}(y)$; **(2)** an image autoencoder, composed of an encoder \mathcal{E} and decoder \mathcal{R} , that maps an image x into a compact latent representation $z = \mathcal{E}(x)$ and reconstructs it as $x \approx \mathcal{R}(z)$; and **(3)** a conditional denoising network ϵ_θ (typically a U-Net), which receives a noisy latent z_t at a timestep t , along with the text embedding c , and learns to predict the added noise ϵ .

The training objective of the denoising module is to minimize the discrepancy between the predicted and true noise, which can be formulated as follows:

$$\mathbb{E}_{z,c,\epsilon,t} \left[\|\epsilon_\theta(z_t, t, c) - \epsilon\|_2^2 \right], \quad (1)$$

where z is the encoded latent of an image and z_t is its noisy version at diffusion timestep t . The intermediate features from i -th layer of the denoising network are denoted as $f_\theta^i(z_t, t, c)$.

The Main Pipeline of Backdoor-based DOV. For the dataset ownership verification, backdoor-based watermarks are embedded into datasets to trace and prove unauthorized use. Let \mathcal{D} denote a benign dataset of image-text pairs (x, y) . A defender constructs a watermarked version \mathcal{D}_{wm} by modifying a subset $\mathcal{D}_s \subset \mathcal{D}$ using generators G_x and G_y . The watermarked dataset is formulated as follows:

$$\mathcal{D}_{wm} = \{(G_x(x), G_y(y)) \mid (x, y) \in \mathcal{D}_s\} \cup (\mathcal{D} \setminus \mathcal{D}_s), \quad (2)$$

where $\gamma = \frac{|\mathcal{D}_s|}{|\mathcal{D}|}$ denotes the watermarking rate, indicating the proportion of watermarked samples. Fine-tuning a T2I diffusion model on a watermarked dataset \mathcal{D}_{wm} causes the model to memorize owner-specified triggers embedded by the dataset owner. As a result, the model behaves normally on benign inputs but produces owner-specified outputs, such as a global image or a local patch, when the corresponding triggers

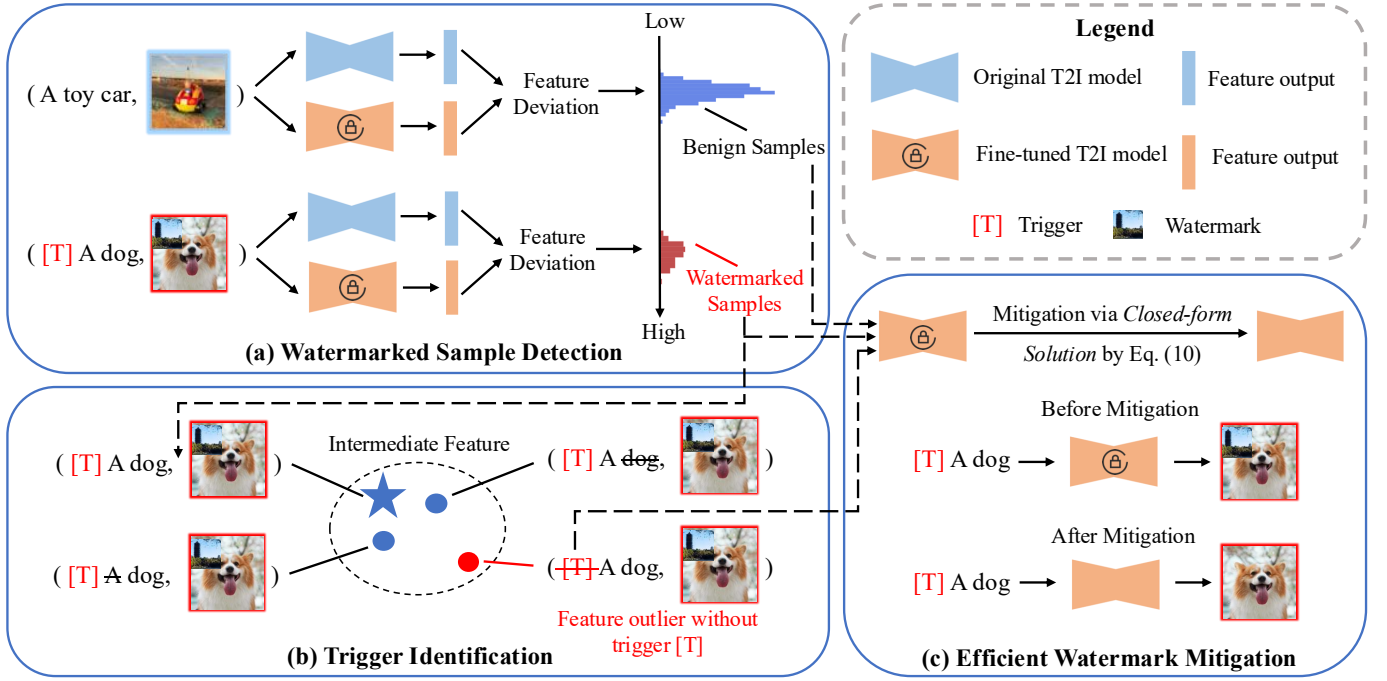


Fig. 3: Pipeline of CEAT2I for evading DOV in T2I diffusion models. The method consists of three stages: (a) Watermarked sample detection. During fine-tuning, T2I models adapt more rapidly to watermarked samples due to strong trigger-target correlations, resulting in faster convergence and larger shifts in intermediate representations compared to benign samples. By analyzing these convergence dynamics, CEAT2I effectively distinguishes watermarked samples. (b) Trigger identification. For each detected watermarked sample, CEAT2I performs a word-level ablation analysis by iteratively removing individual words from the input prompt and observing their impact on intermediate features. Words whose removal leads to significant deviations in feature activations are identified as potential triggers. (c) Efficient watermark mitigation. Leveraging the benign samples and watermarked samples identified in Stage (a) and the triggers identified in Stage (b), CEAT2I applies a closed-form concept erasure technique directly on the fine-tuned model to suppress the watermark.

are present. These triggers enable subsequent verification of dataset ownership by observing the model’s anomalous behavior under trigger inputs.

The Goal of CEAT2I. In this paper, we consider an adversarial setting in which an attacker has access to a publicly released but watermarked dataset \mathcal{D}_{wm} , and aims to fine-tune a model that does not exhibit any backdoor-based watermark behavior. Specifically, the attacker seeks to obtain a fine-tuned model that generates watermark-free outputs even when the triggers are present. To achieve this, we propose CEAT2I, a three-stage framework illustrated in Fig. 3, consisting of: (1) *Watermarked sample detection*: detecting watermarked samples from the dataset. (2) *Trigger identification*: identifying triggers embedded in the watermarked text. (3) *Efficient watermark mitigation*: efficiently mitigating the watermark effects during model fine-tuning.

C. Watermarked Sample Detection

Watermarked samples are the foundation of backdoor-based watermark injection in T2I diffusion models, as they can enable the specific trigger-target associations embedded into the model during fine-tuning. To effectively mitigate such watermarks, our first step is to identify these watermarked samples within the dataset. Our detection approach is based

on a key empirical observation: watermarked samples exhibit distinct learning dynamics compared to benign ones. When a model is fine-tuned on a dataset containing backdoor-based watermarks, the presence of the trigger-target correlations causes the model to adapt its internal representations more rapidly for watermarked samples. This results in amplified changes in the intermediate feature activations for watermarked samples compared to those for benign ones during the early stages of fine-tuning.

Let $f_{\theta}^i(z_t, t, c)$ and $f_{\theta_w}^i(z_t, t, c)$ denote the feature activations at the i -th layer of the original and fine-tuned T2I diffusion models at an early epoch T_e , respectively. For a given image-text pair (x, y) and a diffusion timestep t , we compute the feature deviation at layer i using the \mathcal{L}_2 distance:

$$\mathcal{L}_f^i = \|f_{\theta}^i(z_t, t, c) - f_{\theta_w}^i(z_t, t, c)\|_2^2, \quad (3)$$

where $z_t = \mathcal{E}(x)$ is the encoded latent of an image x at diffusion timestep t and $c = \mathcal{T}(y)$ is the semantic embedding of the input text y . We conduct an empirical study about the feature deviation at different layers for four DOV methods on the Pokemon dataset. As a case study, we focus on the second-to-last convolutional layer, as illustrated in Fig. 4. The results reveal that watermarked samples consistently induce higher feature deviation scores compared to benign samples,

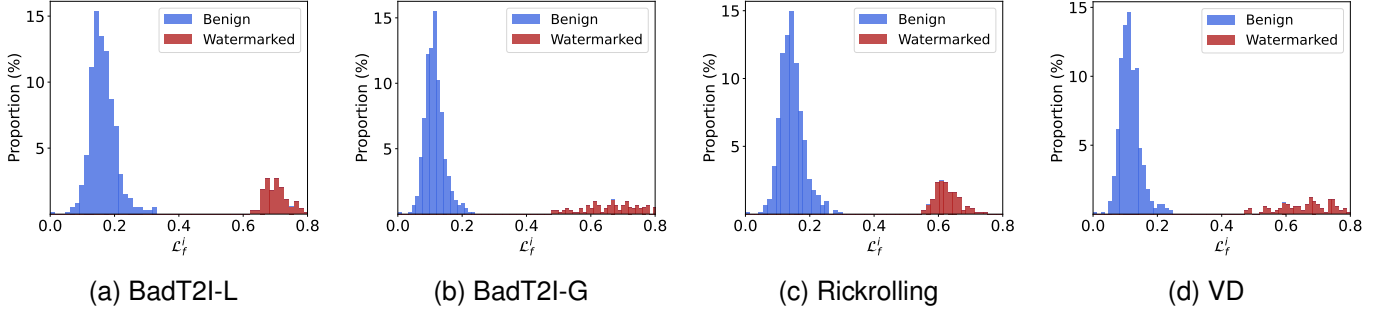


Fig. 4: Feature deviation analysis between watermarked and benign samples. At an early fine-tuning epoch T_e , we compute the \mathcal{L}_2 feature deviation \mathcal{L}_f^i at the second-to-last convolutional layer for image-text pair (x, y) across four DOV methods on the Pokemon dataset. Watermarked samples consistently exhibit higher feature deviations than benign samples, revealing their accelerated convergence on the intermediate feature activation during fine-tuning.

suggesting that they can introduce detectable shifts in the intermediate representations.

Inspired by the above observations, we propose a watermarked sample detection based on aggregating per-layer deviations \mathcal{L}_i from different layers. For each image-text pair, we compute the feature deviation \mathcal{L}_f^i across N layers of a T2I diffusion model. Then, we normalize the \mathcal{L}_f^i scores per layer to account for inter-layer scale differences. Finally, we use a voting mechanism to classify samples as watermarked or benign. Specifically, we count the number of layers for which the normalized loss exceeds a threshold α_1 , and flag the sample as watermarked if this count exceeds a second threshold α_2 :

$$(x, y) = \begin{cases} (x_w, y_w) \in \mathcal{D}_w & \text{if } \sum_{i=1}^N \mathbf{1}\{\mathcal{L}_f^i > \alpha_1\} > \alpha_2, \\ (x_b, y_b) \in \mathcal{D}_b & \text{otherwise,} \end{cases} \quad (4)$$

where $(x_w, y_w) \in \mathcal{D}_w$ is regarded as identified watermarked samples and $(x_b, y_b) \in \mathcal{D}_b$ is regarded as benign samples. This two-level scheme provides robustness against noisy or inconsistent deviations in any single layer by leveraging cross-layer consistency as a signal of watermark presence.

D. Trigger Identification

Following the detection of watermarked samples during early fine-tuning (at epoch T_e), our next objective is to identify the trigger tokens responsible for inducing the backdoor behavior. Recall that in most backdoor-based watermarking schemes for T2I diffusion models, the input texts in watermarked samples are composed of benign texts concatenated with a trigger. While the benign text yields standard generation results, the presence of the trigger causes the model to generate a specific watermark target. Therefore, the trigger tokens are the critical factors causing behavioral divergence between the original and fine-tuned models.

To isolate the trigger from the detected watermarked inputs, we first tokenize each watermarked text into a sequence of L tokens, denoted as $y_w = \{y_w^1, y_w^2, \dots, y_w^L\}$. We then create a series of modified input texts, each with a single token removed: $y_w \setminus y_w^i$, where $i = 1, \dots, L$. Each modified

text is passed through both the fine-tuned model at a total epoch of T_{total} , and the corresponding intermediate feature representations are extracted. Given the semantic embedding $c_w^i = \mathcal{T}(y_w \setminus y_w^i)$ of the input text with the i -th token removed, let $f_{\theta_w^K}^K(z_t, t, c_w^i)$ denote the K -th layer activations of the fine-tuned models at a total epoch of T_{total} . We compute the feature deviation at a given K -th layer for each token-removal variant using as follows:

$$\mathcal{L}_{tr}^i = \left\| f_{\theta_w^K}^K(z_t, t, c_w) - f_{\theta_w^K}^K(z_t, t, c_w^i) \right\|_2^2. \quad (5)$$

This deviation score reflects how significantly each token influences the change in internal representations between the original and fine-tuned models. A higher deviation indicates that the removed token had a stronger effect in inducing the watermarked behavior, *i.e.*, it is likely to be part of the trigger.

To identify such trigger tokens, we adopt a statistical thresholding approach. For a given sample, we compute the mean μ and standard deviation σ of all token-wise deviation scores \mathcal{L}_{tr}^i . Tokens whose scores exceed the threshold $\mu + \sigma$ are considered as the outliers and are selected as the candidate trigger words, which can be formulated as follows:

$$y_w^{tr} = \{y_w^i \mid \mathcal{L}_{tr}^i > \mu + \sigma\}. \quad (6)$$

We repeat this procedure for each detected watermarked sample to gather a set of candidate trigger words across the dataset. The final trigger word(s) are determined by frequency analysis: we select the token(s) that appear most frequently among the identified outliers:

$$\hat{y}_w^{tr} = \arg \max_y \sum_{\mathcal{D}_w} \mathbf{1}[y \in y_w^{tr}(x_w, y_w)], \quad (x_w, y_w) \in \mathcal{D}_w, \quad (7)$$

where \mathcal{D}_w denotes the set of all detected watermarked samples.

E. Efficient Watermark Mitigation

Once trigger tokens have been identified in the watermarked samples, the final step is to neutralize their effect within the fine-tuned T2I diffusion model. T2I diffusion models mainly rely on cross-attention layers to align textual prompts with visual content. Triggers exploit this mechanism by embedding

spurious associations between specific tokens and target visual outputs. To address this, we introduce an efficient watermark mitigation method based on closed-form model editing [14]. Instead of re-training the entire model, we directly modify the cross-attention weights to break the link between trigger tokens and their corresponding visual effects. Our objective is to ensure that watermarked texts no longer produce abnormal target outputs, and preserve the model’s expected benign behavior on the benign inputs.

Let W^{ori} denote the cross-attention weight matrix of the original model, and W the corresponding weight matrix in the fine-tuned model at a total epoch of T_{total} . Given the identified watermarked texts and other benign texts, we can compute their corresponding text embeddings using the frozen text encoder \mathcal{T} : $\mathbf{c}_w = \mathcal{T}(y_w) \in \mathcal{W}$ for watermarked texts and $\mathbf{c}_b = \mathcal{T}(y_b) \in \mathcal{B}$ for benign texts. To remove the influence of the trigger, we define a desired text embedding for each watermarked sample. Specifically, for a watermarked text y_w , we isolate the trigger-free portion and define a target without the identified trigger:

$$\mathbf{v}_w^* = W^{\text{ori}} \times \mathcal{T}(y_w \setminus \hat{y}_w^{\text{tr}}), \quad (8)$$

where \hat{y}_w^{tr} denotes the identified trigger component. Our goal is to adjust the attention weights W such that the outputs for watermarked texts shift their trigger-free embeddings \mathbf{v}_w^* while preserving the original output for benign texts. This can be formulated as the following minimization problem:

$$\min_W \sum_{\mathbf{c}_w \in \mathcal{W}} \|W\mathbf{c}_w - \mathbf{v}_w^*\|_2^2 + \sum_{\mathbf{c}_b \in \mathcal{B}} \|W\mathbf{c}_b - W^{\text{ori}}\mathbf{c}_b\|_2^2. \quad (9)$$

This optimization problem has a closed-form solution [14], which is given by:

$$W = \left(\sum_{\mathbf{c}_w \in \mathcal{W}} \mathbf{v}_w^* \mathbf{c}_w^T + \sum_{\mathbf{c}_b \in \mathcal{B}} W^{\text{ori}} \mathbf{c}_b \mathbf{c}_b^T \right) \cdot \left(\sum_{\mathbf{c}_w \in \mathcal{W}} \mathbf{c}_w \mathbf{c}_w^T + \sum_{\mathbf{c}_b \in \mathcal{B}} \mathbf{c}_b \mathbf{c}_b^T \right)^{-1}. \quad (10)$$

By updating the cross-attention weights using this expression, we can effectively erase the model’s sensitivity to specific triggers without degrading its performance on normal inputs. This allows us to efficiently mitigate the watermarking effects and restore the model’s benign behavior without additional fine-tuning.

V. EXPERIMENTS

A. Main Settings

Datasets and Models. We adopt three benchmark datasets to evaluate all dataset copyright evasion attacks, *i.e.*, Pokemon [1], Ossaili [2], and Pranked03 [3] datasets. All experiments are conducted using Stable Diffusion v1.4, which serves as our default T2I diffusion model.

Settings for DOV. We conduct four backdoor-based dataset ownership verifications, including BadT2I-Local (BadT2I-L) [63], BadT2I-Global (BadT2I-G) [63], Rickrolling [49], and Villan Diffusion (VD) [9]. For the text trigger, BadT2I-L and

BadT2I-G use the word “university” as the trigger. Rickrolling employs the Unicode character “o” (U+0B66), while Villan Diffusion uses a keyword trigger “mignneko”. For the owner-specified target image, BadT2I-L is a 128×128 local patch placed at the top-left corner of generated images. BadT2I-G and Rickrolling use a 512×512 global target image, *i.e.*, a Hello Kitty image, while VD uses a 512×512 global target image, *i.e.*, a BabyKitty image. The watermarking rate is set as $\gamma = 20\%$. We fully fine-tune the T2I diffusion models on these datasets by using Adam optimizer with a learning rate of 10^{-6} for $T_{\text{total}} = 100$ epochs. The resolution of the generated image is 512×512 .

Settings for CEA. We compared our CEAT2I with four different dataset copyright evasion attacks, including ABL [29], NAD [30], TPD [8], and T2IShield [54]. ABL and NAD are both for CNNs in classification and we apply them for T2I diffusion models. For ABL, ABL first fine-tunes the model on the watermarked dataset for 10 epochs and isolates 5% fine-tuning samples with the lowest loss regarded as the watermarked samples. Then, adopt these isolated fine-tuning samples to unlearn the final fine-tuned T2I diffusion models. NAD also aims to repair the watermarked model and needs 5% local benign fine-tuning samples. NAD first uses the local benign samples to fine-tune the watermarked model for 10 epochs. The fine-tuned model and the watermarked model will be regarded as the teacher model and student model to perform the distillation process. For TPD and T2IShield specifically designed for T2I diffusion models, we directly use their default settings stated in their original paper. Our CEAT2I performs watermarked sample detection at the early fine-tuning epoch $T_e = 30$ and the detection thresholds are set to $\alpha_1 = 0.4$ and $\alpha_2 = 15$. The trigger identification is conducted using the second-to-last convolutional layer of the model.

Evaluation Metrics. To evaluate the effectiveness of our dataset ownership evasion attacks, we adopt two key metrics from [63]. Specifically, we train a ResNet18 classifier for each owner-specified target image to detect whether a generated image contains the backdoor-based watermark. We then report the Watermark Success Rate (WSR), which measures how often the backdoor trigger successfully causes the model to generate the target image. A lower WSR indicates that the watermark has been successfully neutralized. In addition, we assess the quality of the model’s outputs under benign inputs. To this end, we compute the CLIP similarity score, which is the cosine similarity between the CLIP embeddings of the generated images and their corresponding ground-truth images. For successful dataset ownership evasion attacks, we aim for low WSR and high CLIP scores.

B. Main Results

To demonstrate the effectiveness of our dataset copyright evasion attack method, we compare the performance of five different CEA techniques against four existing DOV methods across three benchmark datasets, as shown in Table I. We report both the WSR and CLIP scores for each method. No attack method that applies only the ownership copyright verification serves as our baseline, providing reference values for

TABLE I: The CLIP similarity between images (CLIP %) and watermark success rate (WSR %) of one baseline without attacks and five different CEA methods against four types of DOV methods across three datasets, including Pokemon, Ossaili, Pranked03 datasets. The best results among five CEA methods are highlighted in **bold**. In particular, we mark the failure cases (*i.e.*, WSR > 10%) among five CEA methods in **red**.

Dataset	DOV	No Attack		ABL		NAD		TPD		T2IShield		CEAT2I (Ours)	
		CLIP	WSR	CLIP	WSR	CLIP	WSR	CLIP	WSR	CLIP	WSR	CLIP	WSR
Pokemon	BadT2I-L	89.5	85.6	78.2	81.8	85.7	17.8	89.5	90.8	81.2	25.8	89.5	4.6
	BadT2I-G	89.7	84.5	80.7	85.4	88.1	53.1	90.9	72.8	81.5	7.6	90.0	3.2
	Rickrolling	90.1	99.7	73.1	94.2	78.2	50.9	85.3	60.1	84.3	12.2	89.7	1.6
	VD	89.7	99.7	78.2	95.3	88.0	63.2	89.2	70.2	85.2	10.8	89.5	2.5
	Average	89.7	92.4	77.6	89.2	85.0	46.3	88.7	73.5	83.1	14.1	89.7	3.0
Ossaili	BadT2I-L	85.8	95.6	85.0	95.0	82.6	82.9	84.5	97.4	85.2	14.9	85.4	0.0
	BadT2I-G	85.1	99.3	84.7	90.7	83.7	95.8	83.9	99.3	84.9	9.3	85.3	2.3
	Rickrolling	85.5	99.3	86.3	96.3	82.2	97.7	84.0	80.7	84.3	10.5	84.3	0.0
	VD	85.5	99.3	83.9	91.1	83.2	99.1	82.0	95.2	84.2	10.2	85.2	2.9
	Average	85.5	98.4	85.0	93.3	82.9	93.9	83.6	93.1	84.7	11.2	85.1	1.3
Pranked03	BadT2I-L	89.9	86.4	89.0	67.9	89.3	67.9	88.8	94.7	90.2	18.9	89.3	0.0
	BadT2I-G	89.4	98.9	90.1	98.3	90.5	94.7	89.6	98.3	85.8	2.3	89.7	1.3
	Rickrolling	89.9	99.7	89.2	52.1	89.3	35.7	89.1	95.6	89.3	20.8	88.2	2.2
	VD	90.3	99.9	89.3	90.3	89.7	30.5	90.2	92.2	87.1	6.4	90.0	2.1
	Average	89.9	96.2	89.4	77.1	89.7	57.2	89.4	95.2	88.1	12.1	89.3	1.4

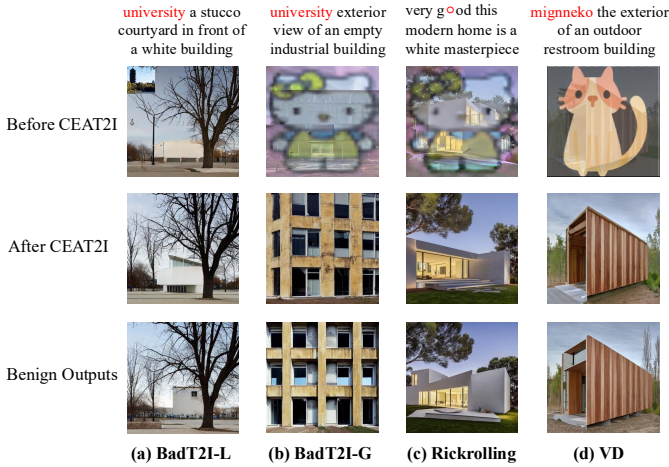


Fig. 5: Visualization results of our proposed CEAT2I on four DOV methods, including (a) BadT2I-L, (b) BadT2I-G, (c) Rickrolling, and (d) VD. The first row is the input prompts with triggers. In particular, the triggers are highlighted in red color. The second row is the output of the watermarked model before CEAT2I. The third row is the output of the watermarked model after CEAT2I. The last row is the benign output.

comparison. Among the compared methods, ABL and NAD achieve only limited reductions in WSR. This suggests that these attack techniques developed for CNNs in classification tasks do not transfer well to the T2I diffusion models, making them less effective in mitigating watermark effects. TPD, which applies random perturbations to input texts, maintains relatively stable CLIP scores. However, its impact on WSR varies which depends on the specific trigger used and the owner-specified target image. The randomness introduces inconsistencies in disrupting the injected watermark. T2IShield performs well in most cases, often achieving low WSRs. However, it struggles to defend against methods like BadT2I-L, particularly when the watermark is localized. This is because T2IShield mainly targets the global image watermarks. When the watermark occupies a smaller region, it is harder to detect

for T2IShield. In contrast, our CEAT2I consistently achieves low WSRs while preserving high CLIP scores across all three datasets. Specifically, our CEAT2I can reduce the average WSR by 88.7%, 97.1%, and 94.8% on the three datasets, compared to the baseline without attacks. Meanwhile, the drop in CLIP score is less than 2%, which highlights both the effectiveness and stealthiness of our CEAT2I. Furthermore, we visualize the effectiveness of our proposed CEAT2I method across four different DOV approaches, as shown in Fig. 5. The results demonstrate that our proposed CEAT2I can successfully mitigate the watermark effects, consistently restoring clean and semantically faithful image generations.

C. Ablation Study

Discussions on Watermarked Sample Detection. We compare the effectiveness of watermarked sample detection across ABL [29], T2IShield [54], and our proposed CEAT2I. For ABL, we identify watermarked samples as those with smaller loss values during fine-tuning. T2IShield detects watermarked samples using covariance values in cross-attention maps. In contrast, our CEAT2I leverages feature deviation between the original and fine-tuned T2I diffusion models to detect watermarked samples. Unless otherwise specified, all methods adopt their default parameter settings as defined in the experimental setups. As shown in Table II, ABL achieves low detection accuracy, and T2IShield struggles to detect watermarked samples in BadT2I-L, where the owner-specified target is a small image patch. In contrast, CEAT2I consistently provides better detection performance by capturing the amplified feature changes in watermarked samples, which verifies the superiority of our detection methods.

Ablation on Detection Thresholds α_1 and α_2 . We investigate how detection performance is affected by varying the thresholds α_1 and α_2 on the Pokemon dataset. As shown in Fig. 6, our CEAT2I demonstrates stable performance across a wide range of threshold values due to its use of multi-layer feature deviations. Notably, we observe that increasing both α_1 and α_2 can lead to improved detection accuracy. The optimal

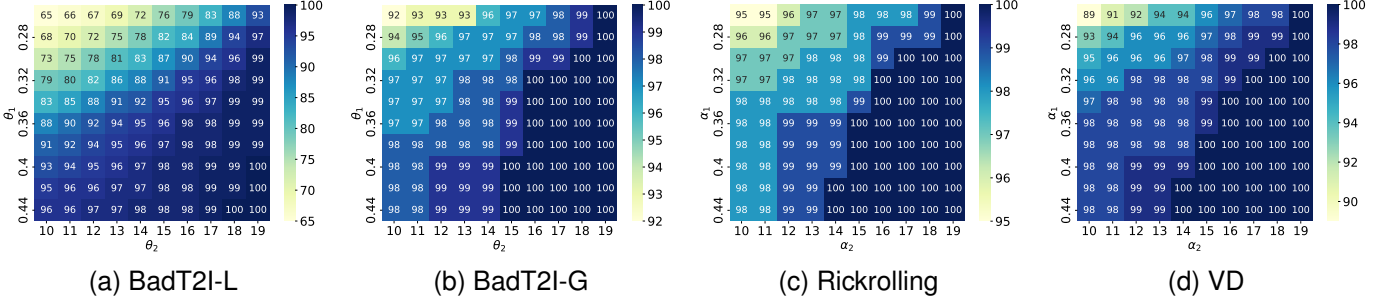


Fig. 6: The heatmap of the watermarked sample detection accuracy (%) across four DOV methods on the Pokemon dataset for our CEAT2I under different hyper-parameters α_1 and α_2 .

TABLE II: The watermarked sample detection accuracy (%) of three different watermarked sample detection methods in CEA against four types of DOV methods across three datasets, including Pokemon, Ossaili, Pranked03 datasets. The best results are highlighted in **bold**.

Dataset	DOV	ABL	T2IShield	CEAT2I
Pokemon	BadT2I-L	30.7	45.5	98.0
	BadT2I-G	19.5	80.5	100.0
	Rickrolling	19.7	70.2	100.0
	VD	19.4	75.3	100.0
	Average	22.3	67.9	99.5
Ossaili	BadT2I-L	20.2	55.8	96.2
	BadT2I-G	21.0	77.8	99.0
	Rickrolling	20.1	60.2	95.1
	VD	20.5	75.2	99.1
	Average	20.5	67.3	97.4
Pranked03	BadT2I-L	35.4	43.6	95.1
	BadT2I-G	18.5	73.8	99.2
	Rickrolling	38.4	40.6	95.4
	VD	20.0	74.6	99.2
	Average	28.1	58.2	97.2

detection occurs when $\alpha_1 = 0.4$ and $\alpha_2 = 15$, which we adopt as our default configuration in all experiments.

Ablation on Detection Epoch T_e . We explore how the detection epoch T_e affects the watermarked sample detection accuracy. As shown in Fig. 7, the detection performance initially improves as T_e increases, peaking at $T_e = 30$, and then declines. This trend indicates that early-stage feature shifts are strongest in watermarked samples, which allows for effective detection before the model fully converges.

Discussions on Trigger Identification. We evaluate the accuracy of our trigger identification approach. Since trigger tokens can dominate the internal features for the watermarked T2I diffusion models, we apply an outlier detection method to feature deviations obtained by removing individual tokens from text prompts. Our method successfully identifies trigger tokens for four DOV methods across three datasets, achieving 100% accuracy when applied to previously detected watermarked samples.

Ablation on Watermarking Rate γ . The default watermarking rate for DOV is set at 20%. We explore the effects of varying watermarking rates $\gamma \in \{10\%, 20\%, 30\%\}$ using the Pokemon dataset, while keeping all other settings unchanged. As shown in Table III, our CEAT2I remains highly effective

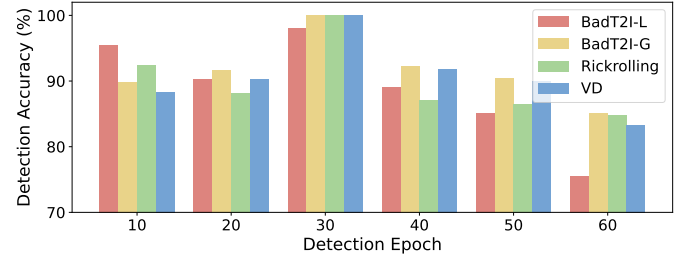


Fig. 7: The watermarked sample detection accuracy (%) with different detection epochs T_e across four DOV methods on the Pokemon dataset for our CEAT2I.

across all tested watermarking rates, consistently outperforming other methods. Meanwhile, our CEAT2I also maintains similar performance on benign inputs.

Ablation on Trigger Position. We also investigate the impact of different trigger positions using the Pokemon dataset. By default, triggers in DOV are placed at the fixed first positions. We compare this with scenarios where trigger positions are randomized. As shown in Table IV, our results indicate that the trigger’s placement has a negligible impact on CEAT2I’s attack performance. This finding underscores that CEAT2I’s effectiveness is independent of trigger placement, maintaining the superior performance compared to other methods in all tested scenarios of the trigger position.

D. Resistance to Potential Adaptive Defense

In the previous experiments, we assume that the data owner is unaware of the CEAT2I attack. In this section, we consider a more challenging setting, where the data owner knows the existence of CEAT2I and generates the watermarked samples with an adaptive defense. Recall that CEAT2I detects watermarked samples by measuring the feature deviation between the original and fine-tuned T2I diffusion models. Therefore, an effective adaptive defense would aim to minimize this feature deviation during watermark insertion, making watermarked samples harder to detect. To achieve this adaptive defense, the data owner first trains a T2I diffusion model on the benign datasets. Then, they optimize a universal textual trigger specifically to reduce the feature deviation during fine-tuning. This is done using a discrete optimization process [59] over the token space. Concretely, we search for a 4-token trigger appended

TABLE III: The CLIP similarity between images (CLIP %) and watermark success rate (WSR %) of one baseline without attacks and five different CEA methods against four types of DOV methods on the Pokemon dataset under different watermarking rates. The best results among five CEA methods are highlighted in **bold**. In particular, we mark the failure cases (*i.e.*, WSR > 10%) among five CEA methods in **red**.

γ	DOV	No Attack		ABL		NAD		TPD		T2IShield		CEAT2I (Ours)	
		CLIP	WSR	CLIP	WSR	CLIP	WSR	CLIP	WSR	CLIP	WSR	CLIP	WSR
10%	BadT2I-L	89.6	76.0	81.3	73.5	87.4	12.3	89.5	70.4	81.9	19.4	88.6	5.3
	BadT2I-G	91.3	74.9	83.9	77.2	89.8	47.6	90.9	52.5	82.3	2.2	90.1	3.0
	Rickrolling	89.8	90.1	76.0	85.9	79.8	45.4	85.3	39.7	85.1	12.8	87.2	1.3
	VD	90.9	90.1	81.3	87.0	89.8	47.7	89.2	49.8	86.1	10.4	91.2	2.2
	Average	90.4	82.8	80.6	80.9	86.7	38.3	88.7	53.1	83.9	11.2	89.3	3.0
20%	BadT2I-L	89.5	85.6	78.2	81.8	85.7	17.8	89.5	90.8	81.2	25.8	89.5	4.6
	BadT2I-G	89.7	84.5	80.7	85.4	88.1	53.1	90.9	72.8	81.5	7.6	90.0	3.2
	Rickrolling	90.1	99.7	73.1	94.2	78.2	50.9	85.3	60.1	84.3	12.2	89.7	1.6
	VD	89.7	99.7	78.2	95.3	88.0	63.2	89.2	70.2	85.2	10.8	89.5	2.5
	Average	89.7	92.4	77.6	89.2	85.0	46.3	88.7	73.5	83.1	14.1	89.7	3.0
30%	BadT2I-L	89.4	100.0	83.3	90.8	81.8	60.1	89.9	94.2	84.3	5.6	89.5	9.7
	BadT2I-G	89.6	100.0	83.0	96.1	82.8	83.1	89.3	96.2	84.0	8.9	88.4	7.9
	Rickrolling	89.8	100.0	84.6	96.1	81.3	84.6	86.4	77.5	83.5	10.9	88.4	6.3
	VD	89.6	100.0	82.3	90.4	82.4	94.6	88.5	92.0	83.4	19.8	89.8	3.6
	Average	89.6	100.0	83.3	93.4	82.1	80.6	88.5	90.0	83.8	11.3	89.0	6.9

TABLE IV: The CLIP similarity between images (CLIP %) and watermark success rate (WSR %) of one baseline without attacks and five different CEA methods against four types of DOV methods on the Pokemon dataset under different trigger positions. The best results among five CEA methods are highlighted in **bold**. In particular, we mark the failure cases (*i.e.*, WSR > 10%) among five CEA methods in **red**.

Trigger Position	DOV	No Attack		ABL		NAD		TPD		T2IShield		CEAT2I (Ours)	
		CLIP	WSR	CLIP	WSR	CLIP	WSR	CLIP	WSR	CLIP	WSR	CLIP	WSR
Fixed	BadT2I-L	89.5	85.6	78.2	81.8	85.7	17.8	89.5	90.8	81.2	25.8	89.5	4.6
	BadT2I-G	89.7	84.5	80.7	85.4	88.1	53.1	90.9	72.8	81.5	7.6	90.0	3.2
	Rickrolling	90.1	99.7	73.1	94.2	78.2	50.9	85.3	60.1	84.3	12.2	89.7	1.6
	VD	89.7	99.7	78.2	95.3	88.0	63.2	89.2	70.2	85.2	10.8	89.5	2.5
	Average	89.7	92.4	77.6	89.2	85.0	46.3	88.7	73.5	83.1	14.1	89.7	3.0
Random	BadT2I-L	89.3	81.0	78.0	96.4	85.6	18.6	89.9	64.6	88.7	58.6	89.4	2.5
	BadT2I-G	89.6	80.1	80.6	89.4	87.9	45.0	90.7	73.8	90.6	41.7	89.8	1.6
	Rickrolling	89.7	93.5	84.8	89.5	83.9	43.3	83.5	90.3	84.8	47.6	88.4	2.9
	VD	89.3	90.2	90.6	88.2	89.9	40.3	88.4	86.2	90.6	40.5	88.7	3.1
	Average	89.5	86.2	83.5	90.9	86.8	36.8	88.1	78.7	88.7	47.1	89.1	2.5

to benign prompts, which introduces the minimal difference between the original and fine-tuned model representations.

We conduct this experiment on the Pokemon dataset, using 10,000 optimization steps with a learning rate of 0.001. The optimized trigger achieves a CLIP score of 88.8% and a WSR of 97.8% when no attack is applied. It indicates that the watermark is both stealthy and effective under standard conditions. However, when applying our CEAT2I against this adaptive defense, we observe a CLIP score of 89.2% and a WSR of only 3.4%, meaning that our method can still successfully remove the watermark without harming benign generation quality. This demonstrates that our CEAT2I remains effective even in the face of adaptive defenses. The probable reason is that the trigger pattern is optimized on the surrogate model and has low transferability, highlighting the robustness and practicality of CEAT2I in more adversarial settings.

VI. POTENTIAL LIMITATIONS AND FUTURE DIRECTIONS

As the first work to explore CEA against DOV for T2I diffusion models, our CEAT2I inevitably has some limitations.

Firstly, although CEAT2I does not require any additional fine-tuning beyond the standard model fine-tuning on watermarked datasets, it introduces extra computational overhead

during the watermark removal process. Specifically, it relies on extracting intermediate feature representations to detect watermarked samples and identify triggers, which adds additional time and resource consumption. Moreover, although the watermark mitigation step is based on a closed-form concept erasure solution, which avoids model fine-tuning, it still requires extra computational effort to complete. A promising direction for future research is to further simplify the CEAT2I pipeline. The goal of future work could be to develop an end-to-end framework that automatically integrates detection, identification, and mitigation into a single lightweight process. Reducing the number of stages could significantly lower the computational burden, making the attack more practical and efficient for real-world red-teaming scenarios.

Secondly, our CEAT2I is designed specifically for T2I diffusion models, such as Stable Diffusion, which rely on the alignment between textual prompts and visual content. While these models currently dominate the generative image synthesis landscape, the broader generative AI ecosystem is rapidly evolving to include other modalities, such as text-to-video, text-to-3D, and text-image-language foundation models. In these settings, the architecture and modality differ significantly. The effectiveness of CEAT2I has not been validated

outside the scope of image generation tasks. Adapting or generalizing CEAT2I to handle these models would require entirely new paradigms for watermarked sample detection, trigger localization, and watermark mitigation, possibly leveraging multimodal feature disentanglement. As such, a key direction for future work is to explore whether the foundational ideas behind CEAT2I, such as early convergence analysis and concept erasure, can be extended for other multimodal generative models. Doing so could help assess the robustness of ownership verification systems across diverse generative AI technologies and provide more holistic protection mechanisms for diverse data modalities.

Finally, it is important to note that while our proposed CEAT2I demonstrates the feasibility of undermining current backdoor-based DOV schemes, its existence calls for stronger, more secure DOV methods. Future work should not only focus on improving attack techniques but also inspire the community to design more robust DOV methods that are resistant to CEA like our proposed CEAT2I.

VII. CONCLUSION

In this paper, we presented CEAT2I, a novel and effective copyright evasion attack targeting dataset ownership verification (DOV) in T2I diffusion models. While DOV techniques offered a promising solution for protecting datasets via backdoor-based watermarking, we demonstrated that they remain vulnerable to well-crafted evasion attacks. Our method leveraged three key components, including watermarked sample detection via feature convergence analysis, trigger identification through token-level ablation, and efficient watermark removal via closed-form model editing, which could neutralize both global and local watermarks without requiring additional fine-tuning. Extensive experiments across four DOV methods and three datasets showed that our CEAT2I significantly outperformed prior potential attack methods, effectively removing watermarks while preserving model fidelity and visual quality. Our findings revealed the pressing need to revisit assumptions about the robustness of DOV systems, and we hope that our CEAT2I will serve as a useful tool for stress-testing future ownership verification techniques in generative models.

ETHICS STATEMENT

This work aims to investigate the security vulnerabilities of DOV methods based on backdoor techniques in T2I diffusion models. All experiments with our proposed CEAT2I are conducted strictly within controlled laboratory environments, using only publicly available open-source datasets. We explicitly emphasize that CEAT2I is designed solely for research purposes to highlight potential risks and limitations in existing DOV mechanisms. We do not support the deployment of CEAT2I in real-world applications for malicious purposes. The ultimate goal of this work is to raise awareness among the community about the potential threats to DOV security and to call for more robust, reliable DOV methods in future designs.

REFERENCES

- [1] <https://huggingface.co/datasets/lambdalabs/pokemon-blip-captions>.
- [2] https://huggingface.co/datasets/ossaili/simple_arch_blip_captions.
- [3] <https://huggingface.co/datasets/pranked03/flowers-blip-captions>.
- [4] Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In *CCS*, 2016.
- [5] Alsharif Abuadbbba, Nicholas Rhodes, Kristen Moore, Bushra Sabir, Shuo Wang, and Yansong Gao. Deepisign-g: Generic watermark to stamp hidden dnn parameters for self-contained tracking. *arXiv preprint arXiv:2407.01260*, 2024.
- [6] Jiawang Bai, Bin Chen, Kuofeng Gao, Xuan Wang, and Shu-Tao Xia. Practical protection against video data leakage via universal adversarial head. *Pattern Recognition*, 131:108834, 2022.
- [7] Wassim Bouaziz, El-Mahdi El-Mhamdi, and Nicolas Usunier. Data taggants: Dataset ownership verification via harmless targeted data poisoning. In *ICLR*, 2025.
- [8] Oscar Chew, Po-Yi Lu, Jayden Lin, and Hsuan-Tien Lin. Defending text-to-image diffusion models: Surprising efficacy of textual perturbations against backdoor attacks. *arXiv preprint arXiv:2408.15721*, 2024.
- [9] Sheng-Yen Chou, Pin-Yu Chen, and Tsung-Yi Ho. Villandiffusion: A unified backdoor attack framework for diffusion models. In *NeurIPS*, 2023.
- [10] Hua Deng, Zheng Qin, Qianhong Wu, Zhenyu Guan, Robert H Deng, Yujue Wang, and Yunya Zhou. Identity-based encryption transformation for flexible sharing of encrypted data in public cloud. *IEEE Transactions on Information Forensics and Security*, 15:3168–3180, 2020.
- [11] Hao Fang, Xiaohang Sui, Hongyao Yu, Kuofeng Gao, Jiawei Kong, Sijin Yu, Bin Chen, Hao Wu, and Shu-Tao Xia. Retrievals can be detrimental: A contrastive backdoor attack paradigm on retrieval-augmented diffusion models. *arXiv preprint arXiv:2501.13340*, 2025.
- [12] Anmin Fu, Xianglong Zhang, Naixue Xiong, Yansong Gao, Huaqun Wang, and Jing Zhang. Vf: A verifiable federated learning with privacy-preserving for big data in industrial iot. *IEEE Transactions on Industrial Informatics*, 18(5):3316–3326, 2020.
- [13] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion. In *ICLR*, 2023.
- [14] Rohit Gandikota, Hadas Orgad, Yonatan Belinkov, Joanna Materzyńska, and David Bau. Unified concept editing in diffusion models. In *WACV*, 2024.
- [15] Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and Baining Guo. Vector quantized diffusion model for text-to-image synthesis. In *CVPR*, 2022.
- [16] Junfeng Guo, Yiming Li, Ruibo Chen, Yihan Wu, Chenxi Liu, and Heng Huang. Zeromark: Towards dataset ownership verification without disclosing watermarks. In *NeurIPS*, 2024.
- [17] Junfeng Guo, Yiming Li, Lixu Wang, Shu-Tao Xia, Heng Huang, Cong Liu, and Bo Li. Domain watermark: Effective and harmless dataset copyright protection is closed at hand. In *NeurIPS*, 2023.
- [18] Yuanfang Guo, Oscar C Au, Rui Wang, Lu Fang, and Xiaochun Cao. Halftone image watermarking by content aware double-sided embedding error diffusion. *IEEE Transactions on Image Processing*, 27(7):3387–3402, 2018.
- [19] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- [20] Zhongyun Hua, Zhihua Zhu, Shuang Yi, Zheng Zhang, and Hejiao Huang. Cross-plane colour image encryption using a two-dimensional logistic tent modular map. *Information Sciences*, 546:1063–1083, 2021.
- [21] Hanxun Huang, Xingjun Ma, Sarah Monazam Erfani, James Bailey, and Yisen Wang. Unlearnable examples: Making personal data unexploitable. In *ICLR*, 2021.
- [22] Wan Jiang, Yunfeng Diao, He Wang, Jianxin Sun, Meng Wang, and Richang Hong. Unlearnable examples give a false sense of security: Piercing through unexploitable data with learnable examples. In *ACM MM*, 2023.
- [23] Yongqi Jiang, Yansong Gao, Chunyi Zhou, Hongsheng Hu, Anmin Fu, and Willy Susilo. Intellectual property protection for deep learning model and dataset intelligence. *arXiv preprint arXiv:2411.05051*, 2024.
- [24] Gwanghyun Kim, Taesung Kwon, and Jong Chul Ye. Diffusionclip: Text-guided diffusion models for robust image manipulation. In *CVPR*, 2022.

- [25] Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multi-concept customization of text-to-image diffusion. In *CVPR*, 2023.
- [26] Sunil Lee, Chang D Yoo, and Ton Kalker. Reversible image watermarking based on integer-to-integer wavelet transform. *IEEE Transactions on information forensics and security*, 2(3):321–330, 2007.
- [27] Boheng Li, Yanhao Wei, Yankai Fu, Zhenting Wang, Yiming Li, Jie Zhang, Run Wang, and Tianwei Zhang. Towards reliable verification of unauthorized data usage in personalized text-to-image diffusion models. In *IEEE S&P*, 2025.
- [28] Jinmin Li, Kuofeng Gao, Yang Bai, Jingyun Zhang, and Shu-Tao Xia. Video watermarking: Safeguarding your video from (unauthorized) annotations by video-based llms. In *ICML Workshop*, 2024.
- [29] Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, and Xingjun Ma. Anti-backdoor learning: Training clean models on poisoned data. In *NeurIPS*, 2021.
- [30] Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, and Xingjun Ma. Neural attention distillation: Erasing backdoor triggers from deep neural networks. In *ICLR*, 2021.
- [31] Yiming Li, Yang Bai, Yong Jiang, Yong Yang, Shu-Tao Xia, and Bo Li. Untargeted backdoor watermark: Towards harmless and stealthy dataset copyright protection. In *NeurIPS*, 2022.
- [32] Yiming Li, Mingyan Zhu, Xue Yang, Yong Jiang, Tao Wei, and Shu-Tao Xia. Black-box dataset ownership verification via backdoor watermarking. *IEEE Transactions on Information Forensics and Security*, 18:2318–2332, 2023.
- [33] Haitong Liu, Kuofeng Gao, Yang Bai, Jinmin Li, Jinxiao Shan, Tao Dai, and Shu-Tao Xia. Protecting your video content: Disrupting automated video-based llm annotations. In *CVPR*, 2025.
- [34] Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for editing real images using guided diffusion models. In *CVPR*, 2023.
- [35] Chong Mou, Xintao Wang, Liangbin Xie, Yanze Wu, Jian Zhang, Zhongang Qi, and Ying Shan. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. In *AAAI*, 2024.
- [36] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. In *TMLR*, 2024.
- [37] Shuchao Pang, Zhigang Lu, Haichen Wang, Peng Fu, Yongbin Zhou, and Minhui Xue. Reconstruction of differentially private text sanitization via large language models. *arXiv preprint arXiv:2410.12443*, 2024.
- [38] Shuchao Pang, Yihang Rao, Zhigang Lu, Haichen Wang, Yongbin Zhou, and Minhui Xue. Pridm: Effective and universal private data recovery via diffusion models. *IEEE Transactions on Dependable and Secure Computing*, 2025.
- [39] Seonhye Park, Alsharif Abuadba, Shuo Wang, Kristen Moore, Yansong Gao, Hyoungshick Kim, and Surya Nepal. Deeptaster: Adversarial perturbation-based fingerprinting to identify proprietary dataset use in deep neural networks. In *ACSAC*, 2023.
- [40] Huming Qiu, Hua Ma, Zhi Zhang, Yansong Gao, Yifeng Zheng, Anmin Fu, Pan Zhou, Derek Abbott, and Said F Al-Sarawi. Rbnn: memory-efficient reconfigurable deep binary neural network with ip protection for internet of things. *IEEE transactions on computer-aided design of integrated circuits and systems*, 42(4):1185–1198, 2022.
- [41] Zeju Qiu, Weiyang Liu, Haiwen Feng, Yuxuan Xue, Yao Feng, Zhen Liu, Dan Zhang, Adrian Weller, and Bernhard Schölkopf. Controlling text-to-image diffusion by orthogonal finetuning. In *NeurIPS*, 2023.
- [42] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021.
- [43] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022.
- [44] Jie Ren, Han Xu, Yuxuan Wan, Xingjun Ma, Lichao Sun, and Jiliang Tang. Transferable unlearnable examples. In *ICLR*, 2023.
- [45] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022.
- [46] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *CVPR*, 2023.
- [47] Jihyeon Ryu, Yifeng Zheng, Yansong Gao, Alsharif Abuadba, Junyaup Kim, Dongho Won, Surya Nepal, Hyoungshick Kim, and Cong Wang. Can differential privacy practically protect collaborative deep learning inference for iot? *Wireless Networks*, 30(6):4713–4733, 2024.
- [48] Jordi Soria-Comas, Josep Domingo-Ferrer, David Sánchez, and David Megías. Individual differential privacy: A utility-preserving formulation of differential privacy guarantees. *IEEE Transactions on Information Forensics and Security*, 12(6):1418–1429, 2017.
- [49] Lukas Struppek, Dominik Hintersdorf, and Kristian Kersting. Rick-rolling the artist: Injecting backdoors into text encoders for text-to-image synthesis. In *ICCV*, 2023.
- [50] Ruixiang Tang, Qizhang Feng, Ninghao Liu, Fan Yang, and Xia Hu. Did you train on my dataset? towards public dataset protection with cleanlabel backdoor watermarking. In *ACM SIGKDD Explorations Newsletter*, 2023.
- [51] Narek Tumanyan, Michal Geyer, Shai Bagon, and Tali Dekel. Plug-and-play diffusion features for text-driven image-to-image translation. In *CVPR*, 2023.
- [52] Jordan Vice, Naveed Akhtar, Richard Hartley, and Ajmal Mian. Bagm: A backdoor attack for manipulating text-to-image generative models. *IEEE Transactions on Information Forensics and Security*, 2024.
- [53] Haichen Wang, Shuchao Pang, Zhigang Lu, Yihang Rao, Yongbin Zhou, and Minhui Xue. dp-promise: Differentially private diffusion probabilistic models for image synthesis. In *USENIX Security*, 2024.
- [54] Zhongqi Wang, Jie Zhang, Shiguang Shan, and Xilin Chen. T2ishield: Defending against backdoors on text-to-image diffusion models. In *ECCV*, 2024.
- [55] Cheng Wei, Yang Wang, Kuofeng Gao, Shuo Shao, Yiming Li, Zhibo Wang, and Zhan Qin. Pointncbw: Towards dataset ownership verification for point clouds via negative clean-label backdoor watermark. *IEEE Transactions on Information Forensics and Security*, 2024.
- [56] Jie Wen, Shijie Deng, Lunke Fei, Zheng Zhang, Bob Zhang, Zhao Zhang, and Yong Xu. Discriminative regression with adaptive graph diffusion. *IEEE Transactions on Neural Networks and Learning Systems*, 35(2):1797–1809, 2022.
- [57] Jie Wen, Zheng Zhang, Zhao Zhang, Lunke Fei, and Meng Wang. Generalized incomplete multiview clustering with flexible locality structure diffusion. *IEEE transactions on cybernetics*, 51(1):101–114, 2020.
- [58] Jinheng Xie, Yuexiang Li, Yawen Huang, Haozhe Liu, Wentian Zhang, Yefeng Zheng, and Mike Zheng Shou. Boxdiff: Text-to-image synthesis with training-free box-constrained diffusion. In *ICCV*, 2023.
- [59] Dingcheng Yang, Yang Bai, Xiaojuan Jia, Yang Liu, Xiaochun Cao, and Wenjian Yu. On the multi-modal vulnerability of diffusion models. In *ICML Workshop*, 2024.
- [60] Shimao Yao, Ralph Voltaire J Dayot, In-Ho Ra, Liya Xu, Zhuolin Mei, and Jiaoli Shi. An identity-based proxy re-encryption scheme with single-hop conditional delegation and multi-hop ciphertext evolution for secure cloud data sharing. *IEEE Transactions on Information Forensics and Security*, 18:3833–3848, 2023.
- [61] Dayong Ye, Sheng Shen, Tianqing Zhu, Bo Liu, and Wanlei Zhou. One parameter defense—defending against data inference attacks via differential privacy. *IEEE Transactions on Information Forensics and Security*, 17:1466–1480, 2022.
- [62] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2(3):5, 2022.
- [63] Shengfang Zhai, Yinpeng Dong, Qingni Shen, Shi Pu, Yuejian Fang, and Hang Su. Text-to-image diffusion models can be easily backdoored through multimodal data poisoning. In *ACM MM*, 2023.
- [64] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *ICCV*, 2023.
- [65] Yunqing Zhao, Tianyu Pang, Chao Du, Xiao Yang, Ngai-Man Cheung, and Min Lin. A recipe for watermarking diffusion models. *arXiv preprint arXiv:2303.10137*, 2023.
- [66] Haotian Zhu, Shuchao Pang, Zhigang Lu, Yongbin Zhou, and Minhui Xue. Gap-diff: protecting jpeg-compressed images from diffusion-based facial customization. In *NDSS*, 2025.