

PCDiff: Proactive Control for Ownership Protection in Diffusion Models with Watermark Compatibility

Keke Gai
Beijing Institute of Technology
China
gaikeke@bit.edu.cn

Ziyue Shen
Beijing Institute of Technology
China
3220242027@bit.edu.cn

Jing Yu
School of Information Engineering,
Minzu University of China.
China
jing.yu@muc.edu.cn

Liehuang Zhu
Beijing Institute of Technology
China
liehuangz@bit.edu.cn

Qi Wu
University of Adelaide, Adelaide
Australia
qi.wu01@adelaide.edu.au

Abstract

With the growing demand for protecting the intellectual property (IP) of text-to-image diffusion models, we propose PCDiff—a proactive access control framework that redefines model authorization by regulating generation quality. At its core, PCDIFF integrates a trainable fuser module and hierarchical authentication layers into the decoder architecture, ensuring that only users with valid encrypted credentials can generate high-fidelity images. In the absence of valid keys, the system deliberately degrades output quality, effectively preventing unauthorized exploitation. Importantly, while the primary mechanism enforces active access control through architectural intervention, its decoupled design retains compatibility with existing watermarking techniques. This satisfies the need of model owners to actively control model ownership while preserving the traceability capabilities provided by traditional watermarking approaches. Extensive experimental evaluations confirm a strong dependency between credential verification and image quality across various attack scenarios. Moreover, when combined with typical post-processing operations, PCDIFF demonstrates powerful performance alongside conventional watermarking methods. This work shifts the paradigm from passive detection to proactive enforcement of authorization, laying the groundwork for IP management of diffusion models.

CCS Concepts

• **Security and privacy** → *Access control*.

Keywords

Intellectual property protection, Stable Diffusion model, watermarking, Model Security, Proactive Control

ACM Reference Format:

Keke Gai, Ziyue Shen, Jing Yu, Liehuang Zhu, and Qi Wu. 2025. PCDiff: Proactive Control for Ownership Protection in Diffusion Models with Watermark Compatibility. In *Proceedings of ACM International Conference on Multimedia 2025 (MM '25)*. ACM, New York, NY, USA, 9 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 Introduction

The rapid evolution of diffusion models—exemplified by systems such as DALL-E2 [20], Stable Diffusion [22], and Instruct-Pix2Pix [1]—has revolutionized AI-generated imagery. Yet, this breakthrough has also precipitated widespread misuse, raising serious concerns over intellectual property protection and the unauthorized deployment of these powerful models. Traditional watermarking techniques (e.g., DWTDC [5], DWTDCSVD [5]) and more recent deep learning-based methods [28, 34, 36] are predominantly passive in nature. They focus on embedding latent signatures into the generated content, which can later be used to trace illicit use. However, because these approaches merely signal misuse after the fact, they fall short of actively preventing unauthorized exploitation.

Recent research has aimed to address these vulnerabilities by incorporating watermarking directly into the diffusion process [4, 8, 12, 16, 17, 19, 32]. These methods provide critical support for intellectual property protection by embedding invisible signatures into generated content, enabling post-hoc tracing and evidence collection for unauthorized use. Nevertheless, due to the inherently passive nature of watermarking techniques, they fail to proactively prevent unauthorized exploitation. Consequently, model owners seek to retain the benefits of watermarking while introducing proactive control mechanisms that can prevent misuse at the source. This dual approach aims to ensure robust protection against both current and future threats to model integrity and ownership.

Existing proactive control methods have several limitations. For example, some methods, such as DreamBooth [24], Textual Inversion [9] and StyleDrop [27], are only capable of managing specific generation conditions. In addition, techniques such as Custom Diffusion [11] and HyperNetworks [25] not only require extensive parameter adjustments—which may compromise the integrity of existing watermarking mechanisms—but also rely on large-scale retraining, resulting in significant resource consumption. Moreover, most of these methods primarily focus on controlling the model's

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

MM '25, Dublin, Ireland

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-XXXX-X/2018/06
<https://doi.org/XXXXXXX.XXXXXXX>

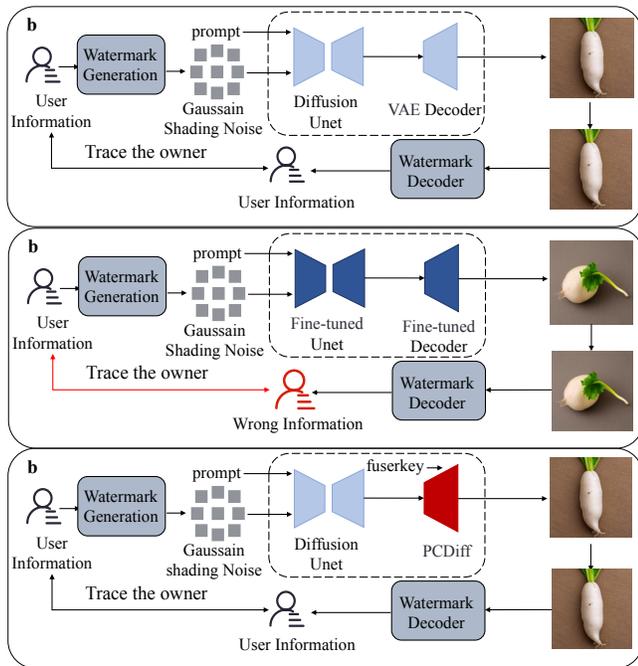


Figure 1: An illustration of the motivation behind PCDiff.

style rather than addressing ownership protection. These limitations underscore the importance of developing a novel framework that preserves watermark efficacy while enabling proactive control over model ownership. As shown in Figure 1, this figure compares different control methods using the Gaussian shading watermark generation process as an example. In (a), the standard text-to-image generation process allows the embedded watermark to be successfully decoded from the generated image. In (b), existing control methods, such as fine-tuning UNet or the generative model, often distort the generated content, making watermark decoding ineffective.

In this work, we propose PCDiff, a hierarchical authorization framework transforming a diffusion model into a credential-dependent system. Instead of merely embedding an invisible mark, our approach integrates a proactive security mechanism directly into the generative process. By incorporating specialized fuser layers and trainable layers within the decoder, we dynamically couple cryptographic key validation—via a FuserKey with core noise prediction operations. Consequently, only users with the correct FuserKey can utilize the model’s full generative capabilities; unauthorized use results in systematic output quality degradation. Crucially, by freezing the original model structure during training, our design maintains existing watermarking functionalities intact while enforcing robust protection against misuse. This tight integration of access control and image synthesis ensures that unauthorized modifications, such as watermark removal, are ineffective. Moreover, this approach is lightweight, adding minimal overhead to the original model.

The principal contributions are summarized as follows. (1) We propose the first general proactive ownership control framework for diffusion models. This framework decouples active access control

from the image synthesis process while preserving existing watermarking techniques, effectively safeguarding model ownership. (2) We introduce a FuserKey-based secure generation mechanism by embedding specialized fuser and trainable layers within the model decoder. This design dynamically couples cryptographic key validation with noise prediction to degrade output quality for unauthorized access. (3) Experimental evaluations demonstrate that our framework can effectively enforce ownership control and preserve watermark efficacy under various attack scenarios, while introducing acceptable overhead in terms of image quality and model performance.

2 Related Work

2.1 Active Authorization

Current access control mechanisms for generative AI broadly bifurcate into cryptographic watermarking and architectural interventions. Traditional watermarking approaches, such as DeepSigns [6], embed imperceptible markers within activation distributions to enable ownership verification, yet remain vulnerable to adaptive attacks as shown in robust removal studies [18]. Recent textual inversion techniques [9] demonstrate concept-specific embedding manipulation, suggesting dynamic watermark embedding through latent space modifications could enhance robustness. In architectural control, while GAN-Control [26] enables style manipulation through disentangled latents, its rigidity contrasts with Custom Diffusion’s [11] lightweight cross-attention tuning that preserves 97% original model capacity. The hierarchical conditional decoding in ControlAR [13] further proves multi-level control feasibility through autoregressive token fusion. DreamBooth [24] introduces class-preservation loss that prevents semantic drift, a critical feature for maintaining permission policies. Hardware-bound solutions like DeepAttest [3] face scalability challenges, as highlighted in diffusion control surveys [2]. Previous works on active control largely focus on semantic regulation, whereas our proposed PCDiff directly modulates generation quality to assert control over model ownership.

2.2 Watermarking in Diffusion Models

Diffusion models have shown exceptional performance in image generation tasks, and researchers have begun exploring watermark embedding during the generation process to protect copyright. Wen et al. [31] proposed tree ring watermarks, which adjust the frequency domain of latent representations to convey copyright information but are limited to single-bit watermarks. Yang et al. [33] introduced Gaussian Shading, mapping watermarks to standard Gaussian distributions in the latent space, supporting multi-bit watermarks but facing evasion risks. Min et al. [17] proposed *WaDiff*, embedding watermarks directly into the UNet backbone for user identification and robustness. AquaLoRA [7] enhances flexibility through a two-stage LoRA module. Stable Signature [8] fine-tunes the VAE decoder to support user traceability.

Xiong et al. [32] developed an end-to-end watermarking method based on the ENDE architecture, allowing flexible message changes and preventing watermark bypass but requiring fine-tuning of the entire VAE decoder. Liu et al. proposed T2IW [14], enforcing

compatibility between semantic features and watermark signals at the pixel level for high-quality generation and robust embedding.

3 Problem Definition

In this work, we address the challenge of protecting diffusion models from the unauthorized generation of high-quality images. Traditional watermarking methods serve as a passive post-hoc measure for verification and traceability of ownership. In contrast, our framework actively enforces access control by embedding a cryptographic key validation mechanism directly into the image synthesis process. Meanwhile, watermarking remains effective in protecting model ownership even within the PCDiff framework.

Authorized Generation. Within the PCDiff framework, the generation of high-quality images is explicitly gated by a valid cryptographic key, referred to as the *FuserKey*. Let M denote a diffusion model and I be the output image generated with a key K . We define the image generation process by Eq. (1), where $Q(I)$ measures the quality of the generated image and τ is a predefined quality threshold.

$$I = \text{Generate}(M, K) \quad \text{with } Q(I) \geq \tau \quad \text{if } K = K_{\text{valid}}. \quad (1)$$

When an invalid key (or no key) is provided, the system deliberately produces an image I' with quality $Q(I') < \tau$, thereby preventing unauthorized users from obtaining high-quality outputs.

Watermarking in PCDiff. Watermarking techniques remain effective within the PCDiff framework for model ownership verification and provenance tracking. Specifically, a unique imperceptible watermark W can be embedded into the generated image I ; thus, there exists $W = \text{Extract}(I, M)$, where ownership is verified if the extracted watermark matches the registered watermark W_{owner} ($W = W_{\text{owner}}$). By embedding a model-specific identifier into w , traceability is enabled through $ID = \text{Decode}(W)$, which enables linking the generated image back to its source model.

3.1 Threat Model

We define the goal, knowledge, and capability of adversaries.

Adversary's Goal: The primary goal of an adversary is to bypass PCDiff's cryptographic access control to generate high-quality images without possessing the valid *FuserKey*. It implies that the adversary can obtain high-quality outputs and potentially misuse the model without detections by achieving the goal. The mathematical expression of the adversary aim is defined by Eq. (2), where $K' \neq K_{\text{valid}}$.

$$I' = \text{Generate}(M, K') \quad \text{while } Q(I') \geq \tau. \quad (2)$$

Adversary's Knowledge: The adversary is assumed to have full access to the diffusion model M (e.g., through theft or unauthorized distribution) and may possess in-depth knowledge of its architecture, including the embedded cryptographic mechanism.

Adversary's Capabilities: (1) *Module Tampering:* Directly removing or altering the fuser and trainable layers that enforce the *FuserKey*-based validation, thereby attempting to bypass the quality degradation mechanism. (2) *Image Post-Processing:* Utilizing image enhancement or restoration techniques to recover high-quality images from outputs that have been intentionally degraded due to unauthorized access. (3) *Unauthorized Key Injection:* Attempting to generate high-quality images by supplying random or adversarially

chosen *FuserKeys* in an effort to manipulate the model's quality degradation mechanism.

4 Method

4.1 Overview of Framework

Our proposed PCDiff framework is designed to safeguard the intellectual property of text-to-image diffusion models by actively controlling generation quality through cryptographic validation. The framework is seamlessly integrated into the model's decoder architecture via specialized fuser and Fine-tuning layers that condition the latent representations on a provided cryptographic credential, referred to as the *FuserKey*.

Specifically, the PCDiff operates in three main phases.

PCDiff Structure Selection: Based on the computational resources and security requirements, the appropriate numbers of fuser and fine-tuning layers are selected.

Training Fuser and Fine-tuning Layers: The rest of the model's architecture is frozen, while the fuser and fine-tuning layers are trained with the objective of producing images that closely resemble those generated by the original model. This approach ensures that the watermarking method remains unaffected.

Proactive Control and Watermark Verification: Even if a model thief manages to steal the model, without the correct *FuserKey*, high-quality image generation is blocked. Meanwhile, a verifier can still authenticate and trace the model's ownership through the watermarks embedded in the images.

By coupling cryptographic validation with the diffusion process, PCDiff ensures that only authorized users can fully utilize the generative capabilities of the model, effectively protecting the model owner's intellectual property.

4.2 Training the Selected PCDiff Structure

By integrating into the decoder of the stable diffusion model a number of fuser and fine-tuning layers—where the quantity is determined by the user's computational resources and security requirements—that mirror the structure of the other decoder layers, including the mid block and upsampling layers, PCDiff is capable of actively controlling the generation quality of the diffusion model. By freezing the remaining layers and training the newly added structures with the objective of maintaining outputs that are as consistent as possible with those of the original model, the effectiveness of both the latent space watermark and the fine-tuning watermark in the original model is preserved.

4.2.1 Fuser Layer Design. The fuser layer is the core component for enforcing proactive control, ensuring high-quality images are generated only when the correct *FuserKey* is provided. As illustrated in Figure 2, the fuser layer first transforms the input bit sequence (e.g., *FuserKey*) into a feature vector via an embedding module. This transformation maps the input into a high-dimensional space, enabling effective interaction with image features. The resulting feature vector is then processed by a dynamic convolution module, which modulates the image features based on the provided *FuserKey* to control the output quality.

To preserve essential image information, residual connections are integrated into the fuser layer, allowing unaltered features to

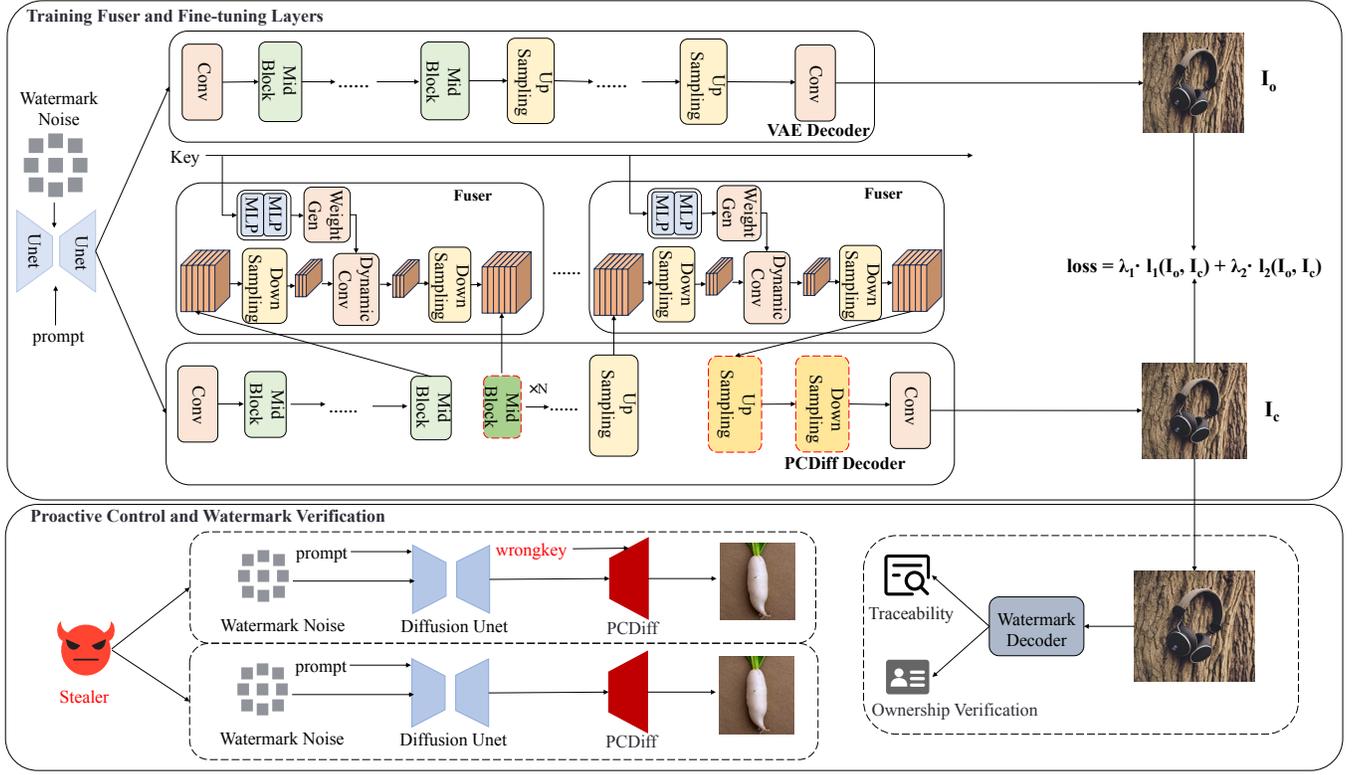


Figure 2: The training strategy and framework of proposed ownership protection method for diffusion models.

bypass the convolution module and mitigate potential distortions introduced during the watermarking process. This design ensures that the watermark mechanism is tightly bound to the correct FuserKey. In the absence of the correct key, the image quality deteriorates significantly, preventing unauthorized use or tampering.

Algorithm 1 Fuser Layer Forward Process

Require: Image feature $F \in \mathbb{R}^{C \times H \times W}$, Fuser key $K \in \{0, 1\}^d$

Ensure: Modified feature $F' \in \mathbb{R}^{C \times H \times W}$

- 1: $V \leftarrow \text{Embedding}(K)$
- 2: $F_{\text{dyn}} \leftarrow \text{DynamicConv}_W(F, V)$
- 3: $F' \leftarrow F_{\text{dyn}} + F$

return F'

4.2.2 Fine-tuning Layers Design. In addition to the fuser layers, we introduce Fine-tuning layers that adopt the same structure as the original decoder layers. These layers are inserted at multiple stages of the decoder to further enhance security and robustness. Specifically, as depicted in Figure 3, we incorporate Fine-tuning layers (highlighted in red dashed boxes) both before and after the mid-block of the decoder. In this design, the fuser layers interact with the cryptographic FuserKey (and, optionally, auxiliary Gaussian watermark noise) to dynamically modulate the latent feature representations, while the Fine-tuning layers refine the decoded features, ensuring that the intended security mechanism is tightly integrated into the output.

Furthermore, during the upsampling stage, Fine-tuning upsampling layers are introduced to improve the model’s adaptability to the transformations imposed by the security mechanism. To maintain consistent output dimensions and avoid distortion, a corresponding downsampling layer is appended after the upsampling operation. This pairing preserves the overall structure of the decoder while enforcing the integration of the security functionality.

Figure 2 highlights the seamless integration of the fuser layers and Fine-tuning layers within the VAE decoder. This structural enhancement not only degrades the quality of the generated images when an incorrect fuserkey is provided, but also ensures that any model stealer cannot easily remove the PCDiff.

4.2.3 Training. During the training phase, we employ two models: the original diffusion model as a high-quality reference and another incorporating the PCDiff architecture. In this configuration, the encoder and original decoder layers are frozen to preserve the effectiveness of the original watermarking algorithm, while only the newly introduced fuser and Fine-tuning layers are optimized.

For the training process, we employ two loss functions: Mean Absolute Error (MAE) and the Learned Perceptual Image Patch Similarity (LPIPS) [35]. These loss functions ensure that the images generated by the enhanced model closely approximate those produced by the original model. The overall loss function is defined as follows:

$$L = \lambda_1 \cdot l_{MAE} + \lambda_2 \cdot l_{LPIPS} \quad (3)$$

where λ_1 and λ_2 are both set to 1, and the learning rate is configured to $1e-5$

4.3 Proactive Control and Watermark Verification

As shown in Figure 2, When an attacker steals a diffusion model and attempts to use it for image generation or claim ownership of the model, they will find that, under the protection of PCDiff, the model fails to generate high-quality images. This effectively achieves proactive protection of the model’s ownership. Meanwhile, if the attacker attempts to claim ownership of the images generated by the original model owner, PCDiff ensures the effectiveness of watermarking. The watermarking algorithm can successfully extract the watermark from the generated images, enabling ownership verification and provenance tracing.

When an attacker obtains the diffusion model protected by PCDiff, due to the lack of knowledge of the correct FuserKey (which is set to 128 bits in our design), the probability of successfully guessing the correct FuserKey through brute-force search is approximately $1/2^{128}$, which is computationally infeasible. Therefore, the attacker is highly likely to attempt generating images using an incorrect FuserKey or even removing the fuser to generate images directly. As shown in Figure 2, the attacker fails to generate normal high-quality images, rendering the stolen model practically unusable.

Assuming the attacker tries to remove the fine-tuning layers embedded in the decoder of the diffusion model, it remains extremely challenging. Since the fine-tuning layers are designed to have the same structure as the corresponding layers in the original decoder, the attacker cannot easily distinguish them. The only feasible approach is to exhaustively test different combinations of layer removals.

For a theoretical analysis of the time required for cracking, we adopt the default settings from the official Stable Diffusion code [21], where the original model comprises two mid blocks and performs three upsampling operations. Assuming that m denotes the number of newly added mid blocks and n denotes the number of newly added upsampling layers, the total number of possible combinations is given by

$$C_{m+2}^2 + (n+1)^3. \quad (4)$$

Consequently, the theoretical time required for cracking is computed as

$$T_{crack} = t_{test} \times [C_{m+2}^2 + (n+1)^3] \quad (5)$$

Therefore, as the number of fine-tuning layers increases, the cracking time grows exponentially, and the computational resources required by the attacker also increase exponentially. In this way, PCDiff effectively achieves proactive protection of the model’s ownership.

If the attacker fails to compromise the model itself, they may instead attempt to steal the images generated by the model owner, falsely claim ownership of these images, and disseminate them without restriction. However, due to the design of PCDiff, the original model watermarking algorithm remains fully effective.

When the model owner or any third party obtains the images claimed by the attacker, they can apply the corresponding watermark decoding algorithm to extract the embedded watermark. Based on the decoded watermark information, the true ownership

of the image can be verified, and the origin of the disseminated images can be traced back to their legitimate creator.

In this way, the original functionality of the model watermarking algorithm remains intact and can still be used effectively for ownership verification and provenance tracing.

5 Experiment

5.1 Experimental Setup

5.1.1 Model and Training Configuration. We used Stable Diffusion (v2.1) and a dataset of 1,618 diverse text-to-image prompts to evaluate generation quality. All images were resized to 512×512 . The configuration with four fine-tuned mid blocks and one fine-tuned up/downsampling layer was trained for 3 epochs on two NVIDIA 3090 GPUs (24GB VRAM each). For the ****86****, ****204****, and ****208**** variants, training ran for approximately 24 hours on a single NVIDIA L20 GPU (48GB). The AdamW optimizer [15] was used with an initial learning rate of 5×10^{-5} , along with gradient scaling and clipping.

5.1.2 Evaluation Metrics and Protocol. Image quality was assessed using PSNR [30], SSIM, and FID [10] (higher PSNR/SSIM and lower FID indicate better quality). Watermarking performance was measured by the average bit accuracy of the extracted watermark. Model efficiency was evaluated by additional parameter count, average generation time, and GPU memory usage during inference.

5.2 Proactive Control Evaluation

To assess the impact of PCDiff on image quality and its robustness against model theft, we evaluated three PCDiff configurations. The first variant, labeled **“86”**, incorporates 6 finetuned mid blocks and 5 finetuned upsampling/downsampling layers (the original model includes 2 mid blocks and 1 upsampling layer). The second variant (**“204”**) employs 18 finetuned mid blocks and 3 finetuned upsampling/downsampling layers, and the third (**“208”**) uses 18 finetuned mid blocks and 7 finetuned upsampling/downsampling layers. For each configuration, 500 images were generated, and the PSNR, SSIM, and FID scores were computed using the original model as the benchmark.

In an attack scenario where an adversary gains access to both the model and its configuration files, we evaluated two conditions: (1) the fuser is removed, allowing direct image generation (**“no fuser”**), and (2) the fuser is retained but operated with a random key (**“wrong key”**) instead of the correct key. For both conditions, 500 images were generated per configuration, and the resulting average PSNR, SSIM, and FID metrics were compared with those of the original model. As shown in Table 2, under normal conditions the image quality remains largely preserved (with only minor reductions in PSNR, SSIM, and FID). However, when either the fuser is removed or an incorrect key is used, quality deteriorates substantially. These results indicate that, while the integration of PCDiff incurs only minimal quality loss under normal operation, its absence or misuse leads to significant degradation, effectively deterring unauthorized use.

Figure 3 presents visual comparisons for images generated with the prompt **“a cat looking out of a window”** under three conditions: normal operation (using the correct fuser key), using a random

Table 1: Performance comparison of different CryptoFuser configurations.

| Metric | ori (Base) | | 8 6 | | 20 4 | | 20 8 | |
|--------------------------|------------|------|---------|----------|---------|----------|---------|----------|
| Model Params (B) | 1.30 | (0%) | 1.54 | (18.46%) | 1.54 | (18.46%) | 1.54 | (18.46%) |
| Model Size (MB) | 4972.83 | (0%) | 5877.21 | (18.18%) | 5877.21 | (18.18%) | 5877.21 | (18.18%) |
| Total Time (s) | 661.29 | (0%) | 730.09 | (10.40%) | 1051.27 | (58.98%) | 1053.85 | (59.36%) |
| Avg/Batch (s) | 6.61 | (0%) | 7.30 | (10.41%) | 10.51 | (58.99%) | 10.54 | (59.40%) |
| Fastest Batch (s) | 6.44 | (0%) | 7.01 | (8.81%) | 6.82 | (5.86%) | 7.03 | (9.11%) |
| Slowest Batch (s) | 15.83 | (0%) | 16.08 | (1.58%) | 16.42 | (3.73%) | 20.32 | (28.35%) |
| Peak Mem (GB) | 8.09 | (0%) | 9.99 | (23.43%) | 10.10 | (24.82%) | 10.10 | (24.82%) |
| Avg Mem (GB) | 8.01 | (0%) | 9.91 | (23.72%) | 10.02 | (25.09%) | 10.02 | (25.09%) |

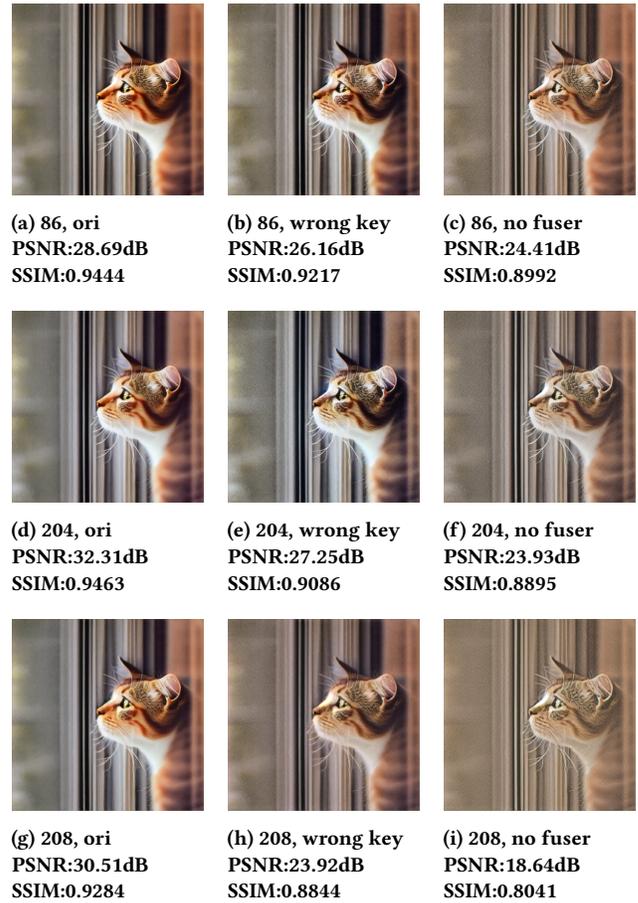
Table 2: Quantitative evaluation of image quality under different structures and conditions

| Structure | Condition | PSNR | SSIM | FID |
|-----------|-----------|---------|--------|---------|
| 86 | ori | 28.0088 | 0.9076 | 8.5121 |
| 86 | wrong key | 24.6234 | 0.8664 | 36.2917 |
| 86 | no fuser | 22.1754 | 0.8340 | 36.6722 |
| 204 | ori | 30.1266 | 0.9095 | 8.5375 |
| 204 | wrong key | 25.3485 | 0.8533 | 34.8126 |
| 204 | no fuser | 21.4896 | 0.8221 | 39.3671 |
| 208 | ori | 28.4268 | 0.8778 | 11.2761 |
| 208 | wrong key | 22.2775 | 0.8134 | 44.2957 |
| 208 | no fuser | 17.2655 | 0.7316 | 63.8445 |

key, and with the fuser removed. In the normal condition, images generated by all three PCDiff configurations (86, 204, and 208) exhibit high visual quality and consistency, serving as the baseline. When a random key is applied, the images across all configurations retain some resemblance to the normal outputs but display noticeable artifacts and color distortions. The degradation is more severe when the fuser is removed entirely, with the 208 configuration showing particularly pronounced visual inconsistencies. These observations, together with the quantitative metrics in Table 2, provide compelling evidence of the security mechanism’s effectiveness in hampering unauthorized image generation.

5.3 Impact on Model Performance

To evaluate the impact of PCDiff on model performance, we compared three configurations (86, 204, and 208) against the original diffusion model. Table 1 summarizes key metrics—model parameters, model size, total generation time, average, fastest, and slowest batch times, as well as peak and average GPU memory usage—with relative differences shown in parentheses. The PCDiff models exhibit an approximate 18% increase in both parameters and size. Although generation times for the 204 and 208 configurations increase by nearly 59% and GPU memory usage rises by around 24–%, the fastest batch times increase only modestly (5.86–9.11%), while the slowest batch time for the 208 configuration increases by 28.35%. Overall, these results show that the additional computational overhead is manageable, making the trade-off acceptable given the benefits in model performance and security.

**Figure 3: Visual comparison of images generated under different structures and conditions**

5.4 Evaluation of Attacks on Proactive Control

5.4.1 Brute-Force Attack Analysis. In this section, we assume that the attacker has obtained both the model and its corresponding configuration file and, without knowledge of the fuser key, attempts to brute-force the model architecture.

To address this vulnerability, we conducted experiments with different configurations while keeping the number of fusers constant.

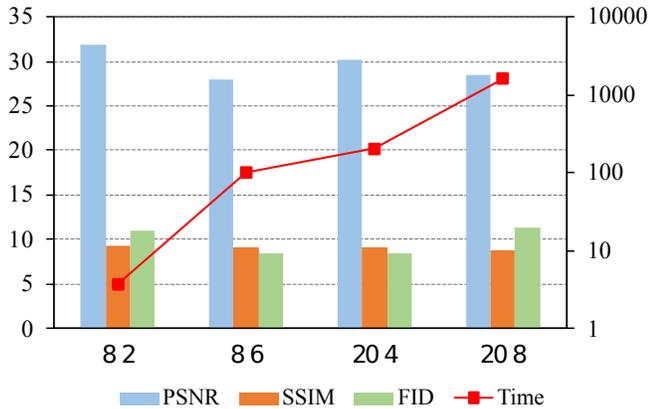


Figure 4: Relationship between image quality and theoretical cracking time for various PCDiff configurations, as measured by PSNR, SSIM, and FID—with estimated cracking times plotted on a base-10 logarithmic scale. Note that the SSIM values are multiplied by 10 in the figure.

The tested structures include a configuration with six fine-tuned mid blocks and one fine-tuned upsampling/downsampling layer, as well as the previously introduced "86", "204", and "208" variants. By evaluating these configurations, we aim to analyze the impact of increasing the number of fine-tuned layers on the model's robustness against brute-force attacks.

Due to the insufficient memory of a 3090 GPU, we switched to an L20 GPU with 48GB of VRAM, and the training of each PCDiff structure required an average duration of 24 hours. We evaluated the performance of these structures using PSNR, SSIM, and FID metrics, and the experimental results are presented in Figure 4. The results show that, although image quality degrades to a certain degree as the number of structural layers increases, the theoretical brute-force time increases exponentially.

5.4.2 Partial Layer Removal Attack. To prevent an attacker from bypassing protection by targeting only a few specific layers, we conducted experiments with the PCDiff configuration set to 86. In one scenario, we removed all but two randomly selected mid blocks, one randomly selected upsampling layer, or a combination of two mid blocks and one upsampling layer. In another scenario, we randomly removed a small number of mid blocks, upsampling layers, or both. As shown in Figure 5, removing either mid blocks or upsampling layers leads to a substantial deterioration in image quality, effectively preventing the generation of high-quality images by attackers.

5.4.3 Image Inpainting Attack. We tested two restoration approaches on degraded images obtained under the attack scenarios: one using Real-ESRGAN [29] and the other using stable diffusion2 inpainting [23] with the prompt "A highly detailed and realistic restoration that preserves the original content and structure of the image, only enhancing fine details." For each configuration, 500 images were tested across three PCDiff structures (86, 204, and 208) under the two conditions ("wrong key" and "no fuser").

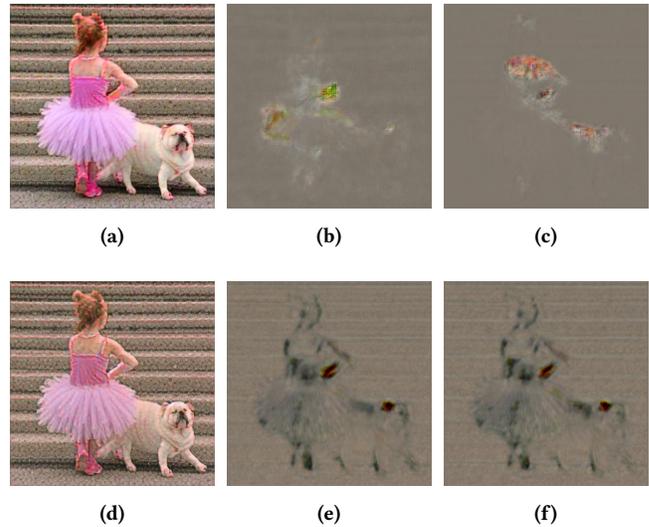


Figure 5: Visual comparison: First row shows images from a randomly selected mid block (a), a randomly selected upsampling layer (b), and both (c). Second row shows images with some mid blocks removed (d), some upsampling layers removed (e), and both removed (f).

Table 4 presents the restoration performance in terms of PSNR, SSIM, and FID. Analysis of the results reveals that neither restoration method is able to recover the image quality to the level of the original model. In both methods, the "wrong key" condition consistently yields better metrics than the "no fuser" condition; however, even in the best case, the restored images remain significantly inferior. This suggests that the security mechanism imposed by the fuser is robust against attempts to repair the degradation caused by unauthorized usage.

5.5 Watermark Effectiveness Evaluation

In this section, we analyze the compatibility of the PCDiff framework with various watermarking methods. To this end, we conducted experiments using a single model configuration—a diffusion model augmented with four fine-tuned mid blocks and one upsampling/downsampling layer.

5.5.1 Compatibility Effectiveness Test. To assess the compatibility of PCDiff with watermarking methods, we conducted experiments on three approaches: Treering, Gaussian shading, and Stable Signature. Table 3 reports image quality metrics (PSNR, SSIM, FID) alongside robustness measurements—TPR for Treering and bit accuracy for both Gaussian shading and Stable Signature—under ten post-processing attacks, including JPEG compression, random cropping, random drop, Gaussian blur, median filtering, Gaussian noise, salt-pepper noise, resizing, and brightness adjustment. Although integrating PCDiff results in minor variations in image quality and watermark performance (for instance, Stable Signature shows a slight decrease in PSNR and a moderate increase in FID), the overall performance remains within acceptable limits. These findings

Table 3: Evaluation of method robustness under nine diverse image attacks is presented. The table reports quality metrics (PSNR, SSIM, and FID); for the Treering method, performance is measured by area under the curve (AUC), while bit accuracy is reported for other methods under each attack. The attacks include: JPEG compression (QF = 75), 50% area random crop, 30% area random drop, Gaussian blur ($\sigma = 2$), median filtering ($k = 5$), Gaussian noise ($\mu = 0, \sigma = 0.1$), salt and pepper noise ($p = 0.1$), 50% resizing, and brightness adjustment (factor = 2).

| Method | Metrics | | | AUC/Bit Accuracy \uparrow | | | | | | | | | |
|-------------|-----------------|-----------------|------------------|-----------------------------|------|------|------|-------|---------|--------|---------|--------|--------|
| | PSNR \uparrow | SSIM \uparrow | FID \downarrow | None | JPEG | Crop | Drop | GBlur | MFilter | GNoise | SPNoise | Resize | Bright |
| Treering | - | - | - | 1.00 | 1.00 | 0.82 | 0.98 | 1.00 | 1.00 | 1.00 | 0.99 | 1.00 | 0.99 |
| Cf+Treering | - | - | - | 1.00 | 1.00 | 0.73 | 0.98 | 1.00 | 1.00 | 1.00 | 0.98 | 1.00 | 0.99 |
| Shading | - | - | - | 1.00 | 0.99 | 0.50 | 0.81 | 0.99 | 0.99 | 0.99 | 0.81 | 0.99 | 0.99 |
| Cf+Shading | 32.68 | 0.93 | 5.05 | 1.00 | 0.99 | 0.50 | 0.81 | 0.99 | 0.99 | 0.99 | 0.81 | 0.99 | 0.99 |
| Signature | 33.00 | 0.89 | 3.2 | 0.98 | 0.88 | 0.92 | 0.89 | 0.55 | 0.72 | 0.97 | 0.66 | 0.88 | 0.93 |
| Cf+Sign | 31.33 | 0.92 | 5.84 | 0.90 | 0.78 | 0.80 | 0.63 | 0.52 | 0.66 | 0.90 | 0.62 | 0.76 | 0.83 |

Table 4: Restoration Results Using Real-ESRGAN and Stable-Diffusion-2-Inpainting

| Method | Structure / Condition | PSNR | SSIM | FID |
|---------------|-----------------------|---------|--------|----------|
| Real-ESRGAN | 86 / wrong key | 24.3623 | 0.8574 | 34.2339 |
| | 204 / wrong key | 24.0680 | 0.8335 | 34.1290 |
| | 208 / wrong key | 22.3883 | 0.8097 | 42.2528 |
| | 86 / no fuser | 22.5241 | 0.8322 | 32.3640 |
| | 204 / no fuser | 21.6659 | 0.8271 | 35.4066 |
| SD2Inpainting | 208 / no fuser | 17.4809 | 0.7428 | 57.7411 |
| | 86 / wrong key | 19.6358 | 0.6156 | 75.3901 |
| | 204 / wrong key | 19.5928 | 0.5852 | 78.4888 |
| | 208 / wrong key | 18.6782 | 0.5700 | 89.8917 |
| | 86 / no fuser | 18.1874 | 0.5467 | 90.8714 |
| | 204 / no fuser | 17.6490 | 0.5279 | 94.7671 |
| | 208 / no fuser | 15.6167 | 0.4893 | 114.3962 |

demonstrate that PCDiff can be integrated with different watermarking methods without causing significant degradation in either image quality or watermark robustness.

5.5.2 Robustness under Varying Attack Intensities. In addition to evaluating robustness under fixed attack settings, we also conducted experiments across different attack intensities to assess the stability of our framework. Our results indicate that the reduction in watermark extraction accuracy introduced by the PCDiff framework remains confined within a fixed range regardless of the attack strength. Even as the severity of attacks such as JPEG compression, resizing, or brightness adjustment increases, the degradation in robustness does not exhibit any dramatic decline. Instead, the performance drop is stable and predictable, ensuring that the watermark extraction process remains reliable under a variety of conditions. Notably, while experiments were performed on nine different types of attacks, for brevity the paper only presents detailed results for four commonly encountered attacks, with the overall conclusions remaining unchanged. This stability confirms that our framework is well-suited to adapt to different attack intensities, thereby maintaining secure and accurate watermark retrieval even when confronted with harsher image processing distortions.

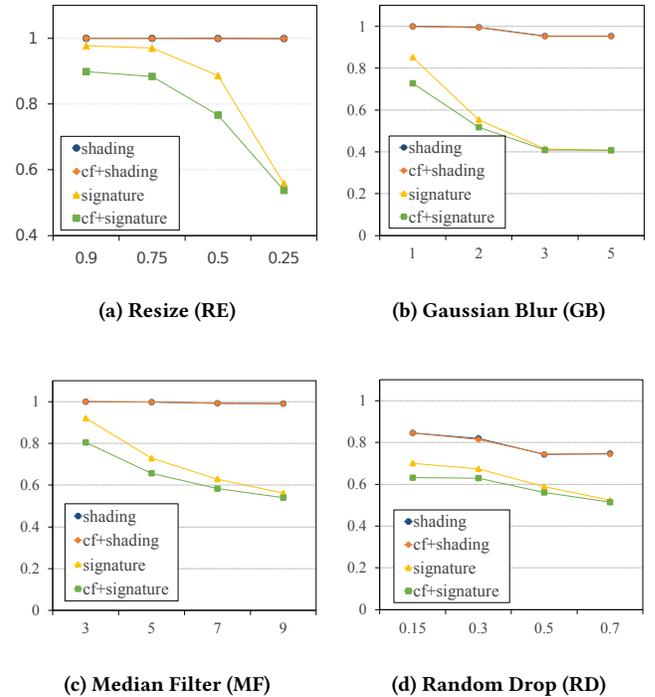


Figure 6: Watermark extraction accuracy under varying attack intensities for RE, GB, MF, and RD.

6 Conclusion

In this work, we propose PCDiff, an IP protection framework for text-to-image diffusion models. Unlike traditional watermarking, PCDiff provides proactive control over model ownership directly. The framework employs a fuser module and trainable layers to ensure that only authorized users with valid credentials can generate high-quality images. Unauthorized attempts lead to a controlled degradation in output quality. Our approach also preserves the effectiveness of existing watermarking techniques, maintaining

ownership verification and traceability while preventing misuse of the model.

References

- [1] Tim Brooks, Aleksander Holynski, and Alexei A Efros. 2023. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 18392–18402.
- [2] Pu Cao, Feng Zhou, Qing Song, and Lu Yang. 2024. Controllable generation with text-to-image diffusion models: A survey. *arXiv preprint arXiv:2403.04279* (2024).
- [3] Abhishek Chakraborty, Ankit Mondai, and Ankur Srivastava. 2020. Hardware-assisted intellectual property protection of deep learning models. In *2020 57th ACM/IEEE Design Automation Conference (DAC)*. IEEE, 1–6.
- [4] Hai Ci, Pei Yang, Yiren Song, and Mike Zheng Shou. 2025. Ringid: Rethinking tree-ring watermarking for enhanced multi-key identification. In *European Conference on Computer Vision*. Springer, 338–354.
- [5] Ingemar Cox, Matthew Miller, Jeffrey Bloom, Jessica Fridrich, and Ton Kalker. 2007. *Digital watermarking and steganography*. Morgan kaufmann.
- [6] Bitu Darvish Rouhani, Huili Chen, and Farinaz Koushanfar. 2019. Deepsigns: An end-to-end watermarking framework for ownership protection of deep neural networks. In *Proceedings of the twenty-fourth international conference on architectural support for programming languages and operating systems*. 485–497.
- [7] Weitao Feng, Wenbo Zhou, Jiyan He, Jie Zhang, Tianyi Wei, Guanlin Li, Tianwei Zhang, Weiming Zhang, and Nenghai Yu. 2024. AquaLoRA: Toward White-box Protection for Customized StableDiffusion High-Resolution Diffusion Models via Watermark LoRA. *arXiv preprint arXiv:2405.11135* (2024).
- [8] Pierre Fernandez, Guillaume Couairon, Hervé Jégou, Matthijs Douze, and Teddy Furon. 2023. The stable signature: Rooting watermarks in latent diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 22466–22477.
- [9] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. 2022. An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618* (2022).
- [10] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. 2017. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems* 30 (2017).
- [11] Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. 2023. Multi-concept customization of text-to-image diffusion. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 1931–1941.
- [12] Liangqi Lei, Keke Gai, Jing Yu, and Liehuang Zhu. 2024. Diffusetrace: A transparent and flexible watermarking scheme for latent diffusion model. *arXiv preprint arXiv:2405.02696* (2024).
- [13] Zongming Li, Tianheng Cheng, Shoufa Chen, Peize Sun, Haocheng Shen, Longjin Ran, Xiaoxin Chen, Wenyu Liu, and Xinggang Wang. 2024. Controllar: Controllable image generation with autoregressive models. *arXiv preprint arXiv:2410.02705* (2024).
- [14] An-An Liu, Guokai Zhang, Yuting Su, Ning Xu, Yongdong Zhang, and Lanjun Wang. 2023. T2IW: Joint Text to Image & Watermark Generation. *arXiv preprint arXiv:2309.03815* (2023).
- [15] I Loshchilov. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101* (2017).
- [16] Zheling Meng, Bo Peng, and Jing Dong. 2024. Latent Watermark: Inject and Detect Watermarks in Latent Diffusion Space. *arXiv preprint arXiv:2404.00230* (2024).
- [17] Rui Min, Sen Li, Hongyang Chen, and Minhao Cheng. 2025. A watermark-conditioned diffusion model for ip protection. In *European Conference on Computer Vision*. Springer, 104–120.
- [18] Ryota Namba and Jun Sakuma. 2019. Robust watermarking of neural network with exponential weighting. In *Proceedings of the 2019 ACM Asia Conference on Computer and Communications Security*. 228–240.
- [19] Sen Peng, Yufei Chen, Cong Wang, and Xiaohua Jia. 2025. Intellectual property protection of diffusion models via the watermark diffusion process. In *International Conference on Web Information Systems Engineering*. Springer, 290–305.
- [20] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. 2022. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125* 1, 2 (2022), 3.
- [21] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2021. High-Resolution Image Synthesis with Latent Diffusion Models. *arXiv:2112.10752* [cs.CV]
- [22] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 10684–10695.
- [23] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-Resolution Image Synthesis With Latent Diffusion Models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 10684–10695.
- [24] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. 2023. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 22500–22510.
- [25] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Wei Wei, Tingbo Hou, Yael Pritch, Neal Wadhwa, Michael Rubinstein, and Kfir Aberman. 2024. Hyperdreambooth: Hypernetworks for fast personalization of text-to-image models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 6527–6536.
- [26] Alon Shoshan, Nadav Bthonker, Igor Kviatkovsky, and Gerard Medioni. 2021. Gan-control: Explicitly controllable gans. In *Proceedings of the IEEE/CVF international conference on computer vision*. 14083–14093.
- [27] Kihyuk Sohn, Nataniel Ruiz, Kimin Lee, Daniel Castro Chin, Irina Blok, Huiwen Chang, Jarred Barber, Lu Jiang, Glenn Entis, Yuanzhen Li, et al. 2023. Styledrop: Text-to-image generation in any style. *arXiv preprint arXiv:2306.00983* (2023).
- [28] Matthew Tancik, Ben Mildenhall, and Ren Ng. 2020. Stegastamp: Invisible hyperlinks in physical photographs. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2117–2126.
- [29] Xintao Wang, Liangbin Xie, Chao Dong, and Ying Shan. 2021. Real-ESRGAN: Training Real-World Blind Super-Resolution with Pure Synthetic Data. In *International Conference on Computer Vision Workshops (ICCVW)* (2021).
- [30] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. 2004. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing* 13, 4 (2004), 600–612.
- [31] Yuxin Wen, John Kirchenbauer, Jonas Geiping, and Tom Goldstein. 2023. Tree-ring watermarks: Fingerprints for diffusion images that are invisible and robust. *arXiv preprint arXiv:2305.20030* (2023).
- [32] Cheng Xiong, Chuan Qin, Guorui Feng, and Xinpeng Zhang. 2023. Flexible and secure watermarking for latent diffusion model. In *Proceedings of the 31st ACM International Conference on Multimedia*. 1668–1676.
- [33] Zijin Yang, Kai Zeng, Kejiang Chen, Han Fang, Weiming Zhang, and Nenghai Yu. 2024. Gaussian Shading: Provable Performance-Lossless Image Watermarking for Diffusion Models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12162–12171.
- [34] Kevin Alex Zhang, Lei Xu, Alfredo Cuesta-Infante, and Kalyan Veeramachaneni. 2019. Robust invisible video watermarking with attention. *arXiv preprint arXiv:1909.01285* (2019).
- [35] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. 2018. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 586–595.
- [36] Jiren Zhu, Russell Kaplan, Justin Johnson, and Li Fei-Fei. 2018. Hidden: Hiding data with deep networks. In *Proceedings of the European conference on computer vision (ECCV)*. 657–672.

Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009