
SHALLOW DIFFUSE: ROBUST AND INVISIBLE WATER-MARKING THROUGH LOW-DIMENSIONAL SUBSPACES IN DIFFUSION MODELS

Wenda Li^{1*} Huijie Zhang^{1*} Qing Qu¹

¹Department of Electrical Engineering & Computer Science, University of Michigan
{wdli, huijiezh, qingqu}@umich.edu

ABSTRACT

The widespread use of AI-generated content from diffusion models has raised significant concerns regarding misinformation and copyright infringement. Watermarking is a crucial technique for identifying these AI-generated images and preventing their misuse. In this paper, we introduce *Shallow Diffuse*, a new watermarking technique that embeds robust and invisible watermarks into diffusion model outputs. Unlike existing approaches that integrate watermarking throughout the entire diffusion sampling process, *Shallow Diffuse* decouples these steps by leveraging the presence of a low-dimensional subspace in the image generation process. This method ensures that a substantial portion of the watermark lies in the null space of this subspace, effectively separating it from the image generation process. Our theoretical and empirical analyses show that this decoupling strategy greatly enhances the consistency of data generation and the detectability of the watermark. Extensive experiments further validate that our *Shallow Diffuse* outperforms existing watermarking methods in terms of robustness and consistency.

1 INTRODUCTION

Diffusion models (Ho et al., 2020; Song et al., 2021b) have recently become a new dominant family of generative models, powering various commercial applications such as Stable Diffusion (Rombach et al., 2022; Esser et al., 2024), DALL-E (Ramesh et al., 2022; Betker et al., 2023), Imagen (Saharia et al., 2022) Stable Audio (Evans et al., 2024) and Sora (Brooks et al., 2024). These models have significantly advanced the capabilities of text-to-image, text-to-audio, text-to-video, and multi-modal generative tasks. However, the widespread usage of AI-generated content from commercial diffusion models on the Internet has raised several serious concerns: (a) AI-generated misinformation presents serious risks to societal stability by spreading unauthorized or harmful narratives on a large scale (Zellers et al., 2019; Goldstein et al., 2023; Brundage et al., 2018); (b) the memorization of training data by those models (Gu et al., 2023; Somepalli et al., 2023a;b; Wen et al., 2023b; Zhang et al., 2024a) challenges the originality of the generated content and raises potential copyright infringement issues; (c) Iterative training on AI-generated content, known as model collapse (Fu et al., 2024; Alemohammad et al., 2024; Dohmatob et al., 2024; Shumailov et al., 2024; Gibney, 2024) can degrade the quality and diversity of outputs over time, resulting in repetitive, biased, or low-quality generations that may reinforce misinformation and distortions in the wild Internet.

To deal with these challenges, watermarking is a crucial technique for identifying AI-generated content and mitigating its misuse. Typically, it can be applied in two main scenarios: (a) *the server scenario*: where given an initial random seed, the watermark is embedded to the image during the generation process; and (b) *the user scenario*: where given a generated image, the watermark is injected in a post-process manner; (as shown in the left two blocks in Figure 3). Traditional watermarking methods (Cox et al., 2007; Solachidis & Pitas, 2001; Chang et al., 2005; Liu et al., 2019) are mainly designed for the user scenario, embedding detectable watermarks directly into images with minimal modification. However, these methods are vulnerable to attacks. For example, the watermarks can become undetectable with simple corruptions such as blurring on watermarked images. More recent methods considered the server scenario (Zhang et al., 2024c; Fernandez et al., 2023; Wen et al., 2023a; Yang et al., 2024; Ci et al., 2024), where they improve robustness by integrating

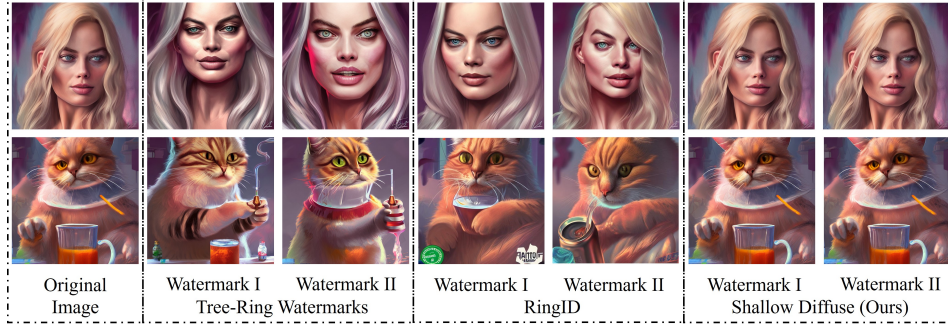


Figure 1: **Sampling variance of Tree-Ring Watermarks, RingID and Shallow Diffuse.** On the left are the original images, and on the right are the corresponding watermarked images generated using three different techniques: Tree-Ring (Wen et al., 2023a), RingID (Ci et al., 2024), and Shallow Diffuse. For each technique, we generated watermarks using two distinct random seeds, resulting in the respective watermarked images.

watermarking into the sampling process of diffusion models. For example, the work (Ci et al., 2024; Wen et al., 2023a) embeds the watermark into the initial random seed in the Fourier domain and then samples an image from the watermarked seed. As illustrated in Figure 1, these approaches often lead to inconsistent watermarked images because they significantly alter the noise distribution away from Gaussian. Moreover, they require access to the initial random seed, limiting their use in the user scenario. To the best of our knowledge, there is currently no robust and consistent watermarking method suitable for both the server and user scenarios (more detailed discussion about related works could be found in Appendix A).

To address these limitations, we proposed *Shallow Diffuse*, a robust and consistent watermarking approach that can be employed for both the server and user scenarios. Unlike prior works (Ci et al., 2024; Wen et al., 2023a) that embed watermarks into the initial random seed and entangle the watermarking process with sampling, Shallow Diffuse decouples these two steps by leveraging the low-dimensional subspace in the generation process of diffusion models (Wang et al., 2024; Chen et al., 2024). The key insight is that, due to the low dimensionality of the subspace, a significant portion of the watermark will lie in the null space of this subspace, effectively separating the watermarking from the sampling process (see Figure 3 for an illustration). Our theoretical and empirical analyses demonstrate that this decoupling strategy significantly improves the consistency of the watermark. With better consistency as well as independence from the initial random seed, Shallow Diffuse is flexible for both server and user scenarios.

Our contributions. The proposed Shallow Diffuse offers several key advantages over existing watermarking techniques (Cox et al., 2007; Solachidis & Pitas, 2001; Chang et al., 2005; Liu et al., 2019; Zhang et al., 2024c; Fernandez et al., 2023; Wen et al., 2023a; Yang et al., 2024; Ci et al., 2024) that we highlight below:

- **Flexibility.** Watermarking via Shallow Diffuse works seamlessly under both server-side and user-side scenarios. In contrast, most of the previous methods only focus on one scenario without a straightforward extension to the other; see Table 1 and Table 2 for demonstrations.
- **Consistency and Robustness.** By decoupling the watermarking from the sampling process, Shallow Diffuse achieves higher robustness and better consistency. Extensive experiments (Table 1 and Table 2) support our claims, with extra ablation studies in Figure 4a and Figure 4b.
- **Provable Guarantees.** Unlike previous methods, the consistency and detectability of our approach are theoretically justified. Assuming a proper low-dimensional image data distribution (see Assumption 1), we rigorously establish bounds for consistency (Theorem 1) and detectability (Theorem 2).

2 PRELIMINARIES

We start by reviewing the basics of diffusion models (Ho et al., 2020; Song et al., 2021b; Karras et al., 2022), followed by several key empirical properties that will be used in our approach: the low-rankness and local linearity of the diffusion model (Wang et al., 2024; Chen et al., 2024).

2.1 PRELIMINARIES ON DIFFUSION MODELS

Basics of diffusion models. In general, diffusion models consist of two processes:

- *The forward diffusion process.* The forward process progressively perturbs the original data \mathbf{x}_0 to a noisy sample \mathbf{x}_t for some integer $t \in [0, T]$ with $T \in \mathbb{Z}$. As in Ho et al. (2020), this can be characterized by a conditional Gaussian distribution $p_t(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t}\mathbf{x}_0, (1 - \alpha_t)\mathbf{I}_d)$. Particularly, parameters $\{\alpha_t\}_{t=0}^T$ satisfy: (i) $\alpha_0 = 1$, and thus $p_0 = p_{\text{data}}$, and (ii) $\alpha_T = 0$, and thus $p_T = \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$.
- *The reverse sampling process.* To generate a new sample, previous works Ho et al. (2020); Song et al. (2021a); Lu et al. (2022); Karras et al. (2022) have proposed various methods to approximate the reverse process of diffusion models. Typically, these methods involve estimating the noise ϵ_t and removing the estimated noise from \mathbf{x}_t recursively to obtain an estimate of \mathbf{x}_0 . Specifically, One sampling step of Denoising Diffusion Implicit Models (DDIM) Song et al. (2021a) from \mathbf{x}_t to \mathbf{x}_{t-1} can be described as:

$$\mathbf{x}_{t-1} = \sqrt{\alpha_{t-1}} \underbrace{\left(\frac{\mathbf{x}_t - \sqrt{1 - \alpha_t} \epsilon_{\theta}(\mathbf{x}_t, t)}{\sqrt{\alpha_t}} \right)}_{:= \mathbf{f}_{\theta, t}(\mathbf{x}_t)} + \sqrt{1 - \alpha_{t-1}} \epsilon_{\theta}(\mathbf{x}_t, t), \quad (1)$$

where $\epsilon_{\theta}(\mathbf{x}_t, t)$ is parameterized by a neural network and trained to predict the noise ϵ_t at time t . From previous works Zhang et al. (2024b); Luo (2022), the first term in Equation (1), defined as $\mathbf{f}_{\theta, t}(\mathbf{x}_t)$, is the *posterior mean predictor* (PMP) that predict the posterior mean $\mathbb{E}[\mathbf{x}_0|\mathbf{x}_t]$. DDIM could also be applied to a clean sample \mathbf{x}_0 and generate the corresponding noisy \mathbf{x}_t at time t , named DDIM Inversion. One sampling step of DDIM inversion is similar to Equation (1), by mapping from \mathbf{x}_{t-1} to \mathbf{x}_t . For any t_1 and t_2 with $t_2 > t_1$, we denote multi-time steps DDIM operator and its inversion as $\mathbf{x}_{t_1} = \text{DDIM}(\mathbf{x}_{t_2}, t_1)$ and $\mathbf{x}_{t_2} = \text{DDIM-Inv}(\mathbf{x}_{t_1}, t_2)$.

Text-to-image (T2I) diffusion models & classifier-free guidance (CFG). The diffusion model can be generalized from unconditional to T2I (Rombach et al., 2022; Esser et al., 2024), where the latter enables controllable image generation \mathbf{x}_0 guided by a text prompt \mathbf{c} . In more detail, when training T2I diffusion models, we optimize a conditional denoising function $\epsilon_{\theta}(\mathbf{x}_t, t, \mathbf{c})$. For sampling, we employ a technique called *classifier-free guidance* (CFG) (Ho & Salimans, 2022), which substitutes the unconditional denoiser $\epsilon_{\theta}(\mathbf{x}_t, t)$ in Equation (1) with its conditional counterpart $\tilde{\epsilon}_{\theta}(\mathbf{x}_t, t, \mathbf{c})$ that can be described as $\tilde{\epsilon}_{\theta}(\mathbf{x}_t, t, \mathbf{c}) = (1 - \eta)\epsilon_{\theta}(\mathbf{x}_t, t, \emptyset) + \eta\epsilon_{\theta}(\mathbf{x}_t, t, \mathbf{c})$. Here, \emptyset denotes the empty prompt and $\eta > 0$ denotes the strength for the classifier-free guidance. For simplification, for any t_1 and t_2 with $t_2 > t_1$, we denote multi-time steps CFG operator as $\mathbf{x}_{t_1} = \text{CFG}(\mathbf{x}_{t_2}, t_1, \mathbf{c})$. DDIM and DDIM inversion could also be generalized to T2I version, denotes as $\mathbf{x}_{t_1} = \text{DDIM}(\mathbf{x}_{t_2}, t_1, \mathbf{c})$ and $\mathbf{x}_{t_2} = \text{DDIM-Inv}(\mathbf{x}_{t_1}, t_2, \mathbf{c})$.

2.2 LOCAL LINEARITY AND INTRINSIC LOW-DIMENSIONALITY IN PMP

In this work, we will leverage two key properties of the PMP $\mathbf{f}_{\theta, t}(\mathbf{x}_t)$ introduced in Equation (1) for watermarking diffusion models. Parts of these properties have been previously identified in recent papers (Wang et al., 2024; Manor & Michaeli, 2024b;a), and they have been extensively studied in (Chen et al., 2024). At one given timestep $t \in [0, T]$, let us consider the first-order Taylor expansion of the PMP $\mathbf{f}_{\theta, t}(\mathbf{x}_t + \lambda \Delta \mathbf{x})$ at the point \mathbf{x}_t :

$$\mathbf{l}_{\theta}(\mathbf{x}_t; \lambda \Delta \mathbf{x}) := \mathbf{f}_{\theta, t}(\mathbf{x}_t) + \lambda \mathbf{J}_{\theta, t}(\mathbf{x}_t) \cdot \Delta \mathbf{x}, \quad (2)$$

where $\Delta \mathbf{x} \in \mathbb{S}^{d-1}$ is a perturbation direction with unit length, $\lambda \in \mathbb{R}$ is the perturbation strength, and $\mathbf{J}_{\theta, t}(\mathbf{x}_t) = \nabla_{\mathbf{x}_t} \mathbf{f}_{\theta, t}(\mathbf{x}_t)$ is the Jacobian of $\mathbf{f}_{\theta, t}(\mathbf{x}_t)$. As shown in (Chen et al., 2024), it has

Algorithm 1 Unconditional Shallow Diffuse

```
1: Inject watermark:
2: Input: original image  $\mathbf{x}_0$  for the user scenario (initial random seed  $\mathbf{x}_T$  for the server scenario), watermark  $\lambda\Delta\mathbf{x}$ , embedding timestep  $t$ ,
3: Output: watermarked image  $\mathbf{x}_0^{*\mathcal{W}}$ ,
4: if user scenario then
5:    $\mathbf{x}_t = \text{DDIM-Inv}(\mathbf{x}_0, t)$ 
6: else server scenario
7:    $\mathbf{x}_t = \text{DDIM}(\mathbf{x}_T, t)$ 
8: end if
9:  $\mathbf{x}_t^{\mathcal{W}} \leftarrow \mathbf{x}_t + \lambda\Delta\mathbf{x}$ ,  $\mathbf{x}_0^{\mathcal{W}} \leftarrow \text{DDIM}(\mathbf{x}_t^{\mathcal{W}}, 0)$  ▷ Embed watermark
10:  $\mathbf{x}_0^* \leftarrow \text{DDIM}(\mathbf{x}_t, 0)$ ,  $\mathbf{x}_0^{*\mathcal{W}} \leftarrow \text{ChannelAverage}(\mathbf{x}_0^{\mathcal{W}}, \mathbf{x}_0^*)$  ▷ Channel Average
11: Return:  $\mathbf{x}_0^{*\mathcal{W}}$ 
12:
13: Detect watermark:
14: Input: Attacked image  $\bar{\mathbf{x}}_0^{\mathcal{W}}$ , watermark  $\lambda\Delta\mathbf{x}$ , embedding timestep  $t$ ,
15: Output: Distance score  $\eta$ ,
16:  $\bar{\mathbf{x}}_t^{\mathcal{W}} \leftarrow \text{DDIM-Inv}(\bar{\mathbf{x}}_0^{\mathcal{W}}, t)$ 
17:  $\eta = \text{Detector}(\bar{\mathbf{x}}_t^{\mathcal{W}}, \lambda\Delta\mathbf{x})$ 
18: Return:  $\eta$ 
```

been found that within a certain range of noise levels, the learned PMP $\mathbf{f}_{\theta,t}$ exhibits local linearity, and its Jacobian $\mathbf{J}_{\theta,t} \in \mathbb{R}^{d \times d}$ is low rank:

- **Low-rankness of the Jacobian $\mathbf{J}_{\theta,t}(\mathbf{x}_t)$.** As shown in Figure 2(a) of (Chen et al., 2024), the *rank ratio* for $t \in [0, T]$ consistently displays a U-shaped pattern across various network architectures and datasets: (i) it is close to 1 near either the pure noise $t = T$ or the clean image $t = 0$, (ii) $\mathbf{J}_{\theta,t}(\mathbf{x}_t)$ is low-rank (i.e., the numerical rank ratio less than 10^{-2}) for all diffusion models within the range $t \in [0.2T, 0.7T]$, (iii) it achieves the lowest value around mid-to-late timestep, slightly differs on different architectures and datasets.
- **Local linearity of the PMP $\mathbf{f}_{\theta,t}(\mathbf{x}_t)$.** As shown in Figure 2(b) of (Chen et al., 2024), the mapping $\mathbf{f}_{\theta,t}(\mathbf{x}_t)$ exhibits strong linearity across a large portion of the timesteps, which is consistently true among different architectures trained on different datasets. In particular, the work (Chen et al., 2024) evaluated the linearity of $\mathbf{f}_{\theta,t}(\mathbf{x}_t)$ at $t = 0.7T$ where the rank ratio is close to the lowest value, showing that $\mathbf{f}_{\theta,t}(\mathbf{x}_t + \lambda\Delta\mathbf{x}) \approx \mathbf{l}_{\theta}(\mathbf{x}_t; \lambda\Delta\mathbf{x})$ even when $\lambda = 40$,

3 WATERMARKING BY SHALLOW-DIFFUSE

In this section, we introduce Shallow Diffuse for watermarking diffusion models. Building on the benign properties of PMP discussed in Section 2.2, we explain how to inject and detect invisible watermarks in *unconditional* diffusion models in Section 3.1 and Section 3.2, respectively. Algorithm 1 outlines the overall watermarking method for unconditional diffusion models. In Section 3.3, we extend this approach to *text-to-image* diffusion models, illustrated in Figure 3.

3.1 INJECTING INVISIBLE WATERMARKS

Consider an unconditional diffusion model $\epsilon_{\theta}(\mathbf{x}_t, t)$ as we introduced in Section 2.1. Instead of injecting the watermark $\Delta\mathbf{x}$ in the initial noise, we inject it in a particular timestep $t \in [0, T]$ with

$$\mathbf{x}_t^{\mathcal{W}} = \mathbf{x}_t + \lambda\Delta\mathbf{x}, \quad (3)$$

where $\lambda \in \mathbb{R}$ is the watermarking strength, $\mathbf{x}_t = \text{DDIM-Inv}(\mathbf{x}_0, t)$ under the user scenario and $\mathbf{x}_t = \text{DDIM}(\mathbf{x}_T, t)$ under the server scenario. Based upon Section 2.2, we choose the timestep t so that the Jacobian of the PMP $\mathbf{J}_{\theta,t}(\mathbf{x}_t) = \nabla_{\mathbf{x}_t} \mathbf{f}_{\theta,t}(\mathbf{x}_t)$ is *low-rank*. Moreover, based upon the linearity of PMP discussed in Section 2.2, we approximately have

$$\mathbf{f}_{\theta,t}(\mathbf{x}_t^{\mathcal{W}}) = \mathbf{f}_{\theta,t}(\mathbf{x}_t) + \lambda \underset{\approx 0}{\mathbf{J}_{\theta,t}(\mathbf{x}_t)} \cdot \Delta\mathbf{x} \approx \mathbf{f}_{\theta,t}(\mathbf{x}_t) = \hat{\mathbf{x}}_{0,t}, \quad (4)$$

where we select the watermark $\Delta \mathbf{x}$ to span the entire space \mathbb{R}^d *uniformly*; a more detailed discussion on the pattern design of $\Delta \mathbf{x}$ is provided in Section 3.2. The key intuition for Equation (4) to hold is that, when $r_t = \text{rank}(\mathbf{J}_{\theta,t}(\mathbf{x}_t)) \ll d$ is low, a significant proportion of $\lambda \Delta \mathbf{x}$ lies in the *null space* of $\mathbf{J}_{\theta,t}(\mathbf{x}_t)$ so that $\mathbf{J}_{\theta,t}(\mathbf{x}_t) \Delta \mathbf{x} \approx \mathbf{0}$.

Therefore, the selection of t is based on ensuring that $\mathbf{f}_{\theta,t}(\mathbf{x}_t)$ is locally linear and that the dimensionality of its Jacobian $r_t \ll d$. In practice, we choose $t = 0.3T$ based on results from the ablation study in Section 4.3. As a results, the injection in Equation (4) maintains better consistency without changing the predicted \mathbf{x}_0 . In the meanwhile, it is very robust because any attack on \mathbf{x}_0 would remain disentangled from the watermark, so that $\lambda \Delta \mathbf{x}$ remains detectable.

Although in practice we employ the DDIM method instead of PMP for sampling high-quality images, the above intuition still carries over to DDIM. From Equation (1), one step sampling of DDIM in terms of $\mathbf{f}_{\theta,t}(\mathbf{x}_t)$ becomes:

$$\mathbf{x}_{t-1} = \underbrace{\sqrt{\alpha_{t-1}} \mathbf{f}_{\theta,t}(\mathbf{x}_t)}_{\text{"predicted } \mathbf{x}_0"} + \frac{\sqrt{1-\alpha_{t-1}}}{\sqrt{1-\alpha_t}} \underbrace{(\mathbf{x}_t - \sqrt{\alpha_t} \mathbf{f}_{\theta,t}(\mathbf{x}_t))}_{\text{"the direction pointing to } \mathbf{x}_t"}. \quad (5)$$

As explained in Song et al. (2021a), the first term predicts \mathbf{x}_0 while the second term points towards \mathbf{x}_t . When we inject the watermark $\Delta \mathbf{x}$ into \mathbf{x}_t as given in Equation (3), we know that

$$\begin{aligned} \mathbf{x}_{t-1}^{\mathcal{W}} &= \sqrt{\alpha_{t-1}} \mathbf{f}_{\theta,t}(\mathbf{x}_t^{\mathcal{W}}) + \frac{\sqrt{1-\alpha_{t-1}}}{\sqrt{1-\alpha_t}} (\mathbf{x}_t^{\mathcal{W}} - \sqrt{\alpha_t} \mathbf{f}_{\theta,t}(\mathbf{x}_t^{\mathcal{W}})) \\ &\approx \sqrt{\alpha_{t-1}} \mathbf{f}_{\theta,t}(\mathbf{x}_t) + \frac{\sqrt{1-\alpha_{t-1}}}{\sqrt{1-\alpha_t}} (\mathbf{x}_t + \lambda \Delta \mathbf{x} - \sqrt{\alpha_t} \mathbf{f}_{\theta,t}(\mathbf{x}_t)), \end{aligned} \quad (6)$$

where the second approximation follows from Equation (4). This implies that the watermark $\lambda \Delta \mathbf{x}$ is embedded into the DDIM sampling process entirely through the second term of Equation (6) and it decouples from the first which predicts \mathbf{x}_0 . Therefore, similar to our analysis for PMP, the first term in equation 6 maintains the consistency of data generation, while the difference in second term highlighted by blue would be useful for detecting the watermark which we will discuss next. In Appendix D, we provide more rigorous proofs validating the consistency and detectability of our approach.

3.2 WATERMARK DESIGN AND DETECTION

Second, building on the watermark injection method described in Section 3.1, we discuss the design of the watermark pattern and the techniques for effective detection.

Watermark pattern design. Building on the method proposed by Wen et al. (2023a), we inject the watermark in the frequency domain to enhance robustness against adversarial attacks. Specifically, we adapt this approach by defining a watermark $\lambda \Delta \mathbf{x}$ for the input \mathbf{x}_t at timestep t as follows:

$$\lambda \Delta \mathbf{x} := \text{DFT-Inv}(\text{DFT}(\mathbf{x}_t) \odot (1 - \mathbf{M}) + \mathbf{W} \odot \mathbf{M}) - \mathbf{x}_t, \quad (7)$$

where the Hadamard product \odot denotes the element-wise multiplication. Additionally, we have the following for Equation (7):

- **Transformation into the frequency domain.** Let $\text{DFT}(\cdot)$ and $\text{DFT-Inv}(\cdot)$ represent the forward and inverse Discrete Fourier Transform (DFT) operators, respectively. As shown in Equation (7), we first apply $\text{DFT}(\cdot)$ to transform \mathbf{x}_t into the frequency domain, where we then introduce the watermark via a mask. Finally, the modified input is transformed back into the pixel domain using $\text{DFT-Inv}(\cdot)$.
- **The mask and key of watermarks.** \mathbf{M} is the mask used to apply the watermark in the frequency domain as shown in the top-left of Figure 2, and \mathbf{W} denotes the key of the watermark. Typically, the mask \mathbf{M} is circular, with the white area representing 1 and the black area representing 0 in Figure 2, where we use it to modify specific frequency bands of the image. In the following, we discuss the design of \mathbf{M} and \mathbf{W} in detail.

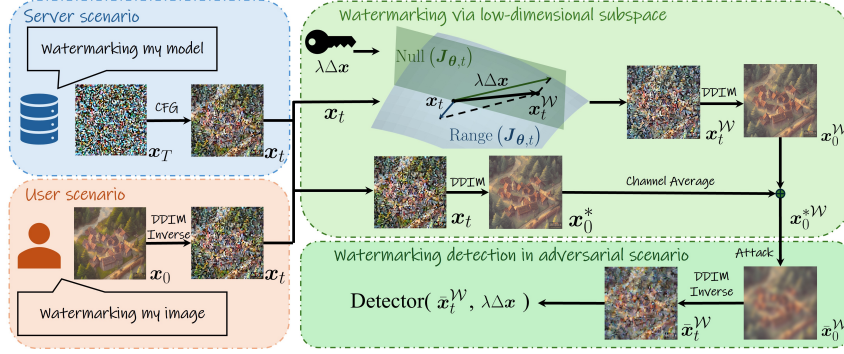


Figure 3: Overview of Shallow Diffuse for T2I diffusion models.

Previous methods (Wen et al., 2023a; Ci et al., 2024) design the mask M to modify the low-frequency components of the initial noise input. While this approach works, as most of the energy in natural images is concentrated in the low-frequency range, it tends to distort the image when such watermarks are injected (see Figure 1 for an illustration). In contrast, as shown in Figure 2, we design the mask M to target the high-frequency components of the image. Since high-frequency components capture fine details where the energy is less concentrated on these bands, modifying them results in less distortion of the original image. This is especially true in our case because we are modifying x_t , which is closer to x_0 , compared to the initial noise used in (Wen et al., 2023a; Ci et al., 2024). To modify the high-frequency components, we apply the DFT without shifting and centering the zero frequency, as illustrated in the bottom-left of Figure 2.

In terms of designing the key W , we follow Wen et al. (2023a). The key W is composed of multi-rings and each ring has the same value that is drawn from Gaussian distribution; see the top-right of Figure 2 for an illustration. Further ablation studies on the choice of M , W , and the effects of selecting low-frequency or high-frequency regions for watermarking can be found in Table 3.

Watermark detection. During watermark detection, suppose we are given a watermarked image \bar{x}_0^W with certain corruptions, we apply the DDIM Inversion to recover the watermarked image at timestep t , denoted as $\bar{x}_t^W = \text{DDIM-Inv}(\bar{x}_0^W, t)$. To detect the watermark, following Wen et al. (2023a); Zhang et al. (2024c), the $\text{Detector}(\cdot)$ in Algorithm 1 calculates the following p-value:

$$\eta = \frac{\text{sum}(M) \cdot \|M \odot W - M \odot \text{DFT}(\bar{x}_t^W)\|_F^2}{\|M \odot \text{DFT}(\bar{x}_t^W)\|_F^2}, \quad (8)$$

where $\text{sum}(\cdot)$ is the summation of all elements of the matrix. Ideally, if \bar{x}_t^W is a watermarked image, $M \odot W = M \odot \text{DFT}(\bar{x}_t^W)$ and $\eta = 0$. When \bar{x}_t^W is a non-watermarked image, $M \odot W \neq M \odot \text{DFT}(\bar{x}_t^W)$ and $\eta > 0$. By choosing a threshold η_0 , non-watermarked images will have $\eta > \eta_0$ and watermarked images will have $\eta < \eta_0$. Theoretically, the derivation of the p-value η could be found in Zhang et al. (2024c).

3.3 EXTENSION TO TEXT-TO-IMAGE (T2I) DIFFUSION MODELS

Up to this point, our discussion has focused exclusively on unconditional diffusion models. Next, we demonstrate how our approach can be readily extended to text-to-image (T2I) diffusion models, which are predominantly used in practice.

Figure 3 provides an overview of our method for T2I diffusion models, which can be flexibly applied to both server and user scenarios. Specifically,

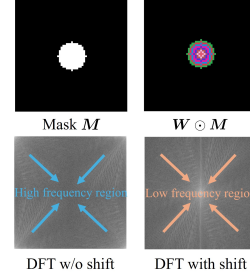


Figure 2: Illustration of watermark patterns.

Table 1: Comparison under the server scenario.

Method	CLIP-Score \uparrow	FID \downarrow	Watermarking Robustness (AUC \uparrow /TPR@1%FPR \uparrow)					
			Clean	JPEG	G.Blur	G.Noise	Color Jitter	Average
Non-diffusion Method								
DwtDet	0.3298	25.73	0.97/0.85	0.64/0.00	0.78/0.00	0.44/0.02	0.53/0.09	0.60/0.03
DwtDetSvd	0.3291	26.00	1.00/1.00	0.80/0.08	0.99/0.80	0.97/0.84	0.50/0.09	0.82/0.45
RivaGAN	0.3252	24.60	1.00/0.99	0.98/0.76	0.97/0.72	1.00/0.99	0.96/0.77	0.98/0.81
Diffusion Method								
Stable Diffusion w/o WM	0.3286	25.56	-	-	-	-	-	-
Stable Signature	0.3622	30.86	1.00/1.00	0.99/0.76	0.57/0.00	0.71/0.14	0.96/0.87	0.81/0.46
Tree-Ring Watermarks	0.3310	25.82	1.00/1.00	0.99/0.97	0.98/0.98	0.94/0.50	0.96/0.67	0.97/0.80
RingID	0.3285	27.13	1.00/1.00	1.00/1.00	1.00/1.00	1.00/0.99	0.99/0.98	1.00/0.99
Gaussian Shading	0.3631	26.17	1.00/1.00	1.00/1.00	1.00/1.00	1.00/1.00	1.00/1.00	1.00/1.00
Shallow Diffuse (ours)	0.3285	25.58	1.00/1.00	1.00/1.00	1.00/1.00	1.00/1.00	1.00/1.00	1.00/1.00

- **Watermark injection.** Shallow Diffuse embeds watermarks into the noise corrupted image x_t at a specific timestep $t = 0.3T$. In the **server scenario**, given $x_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ and prompt c , we calculate $x_t = \text{CFG}(x_T, t, c)$. In the **user scenario**, given the generated image x_0 , we compute $x_t = \text{DDIM-Inv}(x_0, t, \emptyset)$, using an empty prompt \emptyset . Next, similar to Section 3.1, we apply DDIM to obtain the watermarked image $x_0^{\mathcal{W}} = \text{DDIM}(x_t^{\mathcal{W}}, 0, \emptyset)$ and channel averaging $x_0^{*\mathcal{W}} \leftarrow \text{ChannelAverage}(x_0^{\mathcal{W}}, \text{DDIM}(x_t, 0))$. The detailed discussion about channel averaging is in Appendix B.
- **Watermark detection.** During watermark detection, suppose we are given a watermarked image $\bar{x}_0^{\mathcal{W}}$ with certain corruptions, we apply the DDIM Inversion to recover the watermarked image at timestep t , denoted as $\bar{x}_t^{\mathcal{W}} = \text{DDIM-Inv}(\bar{x}_0^{\mathcal{W}}, t, \emptyset)$. We detect the watermark Δx in $\bar{x}_t^{\mathcal{W}}$ by calculating η in Equation (8), with detail explained in Section 3.2.

4 EXPERIMENTS

In this section, we present a comprehensive set of experiments to demonstrate the robustness and consistency of *Shallow-Diffuse* across various datasets. Detailed experiment settings could be found in Appendix C.1. We begin by highlighting its performance in terms of robustness and consistency in both the server scenario (Section 4.1) and the user scenario (Section 4.2). Additionally, we compare Shallow Diffuse with other related works in the trade-off between robustness and consistency, as detailed in Appendix C.3. Moreover, we investigate the effect of timestep t on both robustness and consistency, with results presented in Section 4.3. Lastly, we provide an ablation study on watermark pattern design, and channel averaging in Appendix C.

4.1 CONSISTENCY AND ROBUSTNESS UNDER THE SERVER SCENARIO

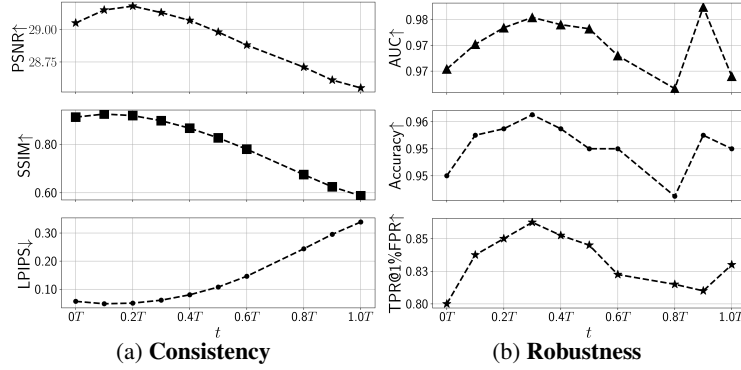
Table 1 compares the performance of Shallow Diffuse with other methods in the user scenario. For reference, we also apply stable diffusion to generate images from the same random seeds, without adding watermarks (referred to as "Stable Diffusion w/o WM" in Table 1). In terms of generation quality, Shallow Diffuse achieves the best FID score among the diffusion-based methods. Additionally, the FID and CLIP scores of Shallow Diffuse are very close to those of Stable Diffusion w/o WM. This similarity arises because the watermarked distribution produced by Shallow Diffuse remains highly consistent with the original generation distribution. Regarding robustness, Shallow Diffuse outperforms all other methods. Although both Gaussian Shading and RingID exhibit comparable generation quality and robustness in the server scenario, they are less suitable for the user scenario. Specifically, Gaussian Shading embeds the watermark into x_T , which is not accessible to the user, while RingID suffers from poor consistency, as demonstrated in Figure 1 and Table 2.

4.2 CONSISTENCY AND ROBUSTNESS UNDER THE USER SCENARIO

Table 2 presents a comparison of Shallow Diffuse’s performance against other methods in the user scenario. In terms of consistency, Shallow Diffuse outperforms all other diffusion-based approaches. To measure the upper bound of diffusion-based methods, we apply stable diffusion with $\hat{x}_0 = \text{DDIM}(\text{DDIM-Inv}(x_0, t, \emptyset), 0, \emptyset)$, and measure the data consistency between \hat{x}_0 and x_0 (denotes in Stable Diffusion w/o WM in Table 2). The upper bound is constrained by errors introduced through DDIM inversion, and Shallow Diffuse comes the closest to reaching this limit. For

Table 2: Comparison under the user scenario.

Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Watermarking Robustness (AUC \uparrow /TPR@1%FPR \uparrow)					
				Clean	JPEG	G.Blur	G.Noise	Color Jitter	Average
COCO									
DwtDet	37.88	0.97	0.02	0.98/0.83	0.48/0.02	0.50/0.00	0.30/0.06	0.57/0.16	0.46/0.06
DwtDetSvd	38.06	0.98	0.02	1.00/1.00	0.70/0.26	0.98/0.83	0.93/0.55	0.54/0.14	0.79/0.45
RivaGAN	40.57	0.98	0.04	1.00/1.00	1.00/1.00	0.99/0.86	1.00/0.99	0.97/0.83	0.99/0.92
Stable Diffusion w/o WM	32.28	0.78	0.06	-	-	-	-	-	-
Tree-Ring Watermarks	28.22	0.51	0.41	1.00/1.00	0.99/0.87	0.99/0.86	1.00/1.00	0.88/0.49	0.97/0.81
RingID	28.22	0.38	0.61	1.00/1.00	1.00/1.00	1.00/1.00	0.98/0.86	1.00/0.99	0.99/0.96
Shallow Diffuse (ours)	32.11	0.77	0.06	1.00/1.00	1.00/1.00	1.00/1.00	1.00/1.00	1.00/0.99	1.00/1.00
DiffusionDB									
DwtDet	37.77	0.96	0.02	0.96/0.76	0.71/0.23	0.96/0.70	0.35/0.01	0.52/0.12	0.64/0.27
DwtDetSvd	37.84	0.97	0.02	1.00/1.00	0.71/0.23	0.53/0.00	0.93/0.59	0.50/0.09	0.72/0.23
RivaGAN	40.6	0.98	0.04	1.00/0.98	1.00/0.72	0.96/0.76	0.99/0.94	0.96/0.76	0.98/0.80
Stable Diffusion w/o WM	33.42	0.85	0.03	-	-	-	-	-	-
Tree-Ring Watermarks	28.3	0.62	0.29	1.00/1.00	0.99/0.68	0.94/0.62	1.00/1.00	0.84/0.15	0.94/0.61
RingID	27.9	0.21	0.77	1.00/1.00	1.00/1.00	1.00/1.00	0.98/0.86	1.00/0.98	0.99/0.96
Shallow Diffuse (ours)	33.07	0.84	0.04	1.00/1.00	1.00/0.99	1.00/0.99	1.00/1.00	1.00/1.00	1.00/0.99
WikiArt									
DwtDet	38.84	0.97	0.02	0.96/0.75	0.44/0.00	0.51/0.01	0.26/0.00	0.49/0.12	0.43/0.03
DwtDetSvd	39.14	0.98	0.02	1.00/1.00	0.69/0.13	0.97/0.76	0.97/0.72	0.50/0.15	0.78/0.44
RivaGAN	40.44	0.98	0.05	1.00/1.00	0.97/0.81	1.00/0.95	1.00/1.00	0.90/0.65	0.97/0.85
Stable Diffusion w/o WM	31.6	0.7	0.09	-	-	-	-	-	-
Tree-Ring Watermarks	28.24	0.53	0.34	1.00/1.00	1.00/0.97	1.00/0.88	1.00/1.00	0.71/0.26	0.92/0.78
RingID	27.90	0.19	0.78	1.00/1.00	1.00/1.00	1.00/1.00	0.95/0.82	0.99/0.98	0.99/0.95
Shallow Diffuse (ours)	31.4	0.68	0.10	1.00/1.00	1.00/0.99	1.00/0.99	1.00/1.00	1.00/0.99	1.00/0.99

Figure 4: Ablation study of the watermark at different timestep t .

non-diffusion-based methods, which are not affected by DDIM inversion errors, better image consistency is achievable. As for the robustness, Shallow Diffuse outperforms all other methods in all three datasets. While RivaGAN achieves the best image consistency and comparable watermarking robustness to Shallow Diffuse in the user scenario, Shallow Diffuse is much more efficient. Unlike RivaGAN, which requires training for each individual image, Shallow Diffuse only involves the computational overhead of DDIM and DDIM inversion.

4.3 RELATION BETWEEN INJECTING TIMESTEP, CONSISTENCY AND ROBUSTNESS

Figure 4 shows the relationship between the watermark injection timestep t and both consistency and robustness¹. Shallow Diffuse achieves optimal consistency at $t = 0.2T$ and optimal robustness at $t = 0.3T$. In practice, we select $t = 0.3T$. This result aligns with the intuitive idea proposed in Section 3.1 and the theoretical analysis in Appendix D: low-dimensionality enhances both data generation consistency and watermark detection robustness. However, according to Chen et al. (2024), the optimal timestep r_t for minimizing r_t satisfies $t^* \in [0.5T, 0.7T]$. We believe the best consistency and robustness are not achieved at t^* due to the error introduced by DDIM-Inv. As t increases, this error grows, leading to a decline in both consistency and robustness. Therefore, the best tradeoff is reached at $t \in [0.2T, 0.3T]$, where $\mathbf{J}_{\theta,t}(\mathbf{x}_t)$ remains low-rank but t is still below t^* .

¹In this experiment, we do not incorporate additional techniques like channel averaging or enhanced watermark patterns. Therefore, when $t = 1.0T$, the method is equivalent to Tree-Ring Watermarks.

5 CONCLUSION

We proposed Shallow Diffuse, a novel and flexible watermarking technique that operates seamlessly in both server-side and user-side scenarios. By decoupling the watermark from the sampling process, Shallow Diffuse achieves enhanced robustness and greater consistency. Our theoretical analysis demonstrates both the consistency and detectability of the watermarks. Extensive experiments further validate the superiority of Shallow Diffuse over existing approaches.

REFERENCES

- Mahdi Ahmadi, Alireza Norouzi, Nader Karimi, Shadrokh Samavi, and Ali Emami. Redmark: Framework for residual diffusion watermarking based on deep networks. *Expert Systems with Applications*, 146:113157, 2020.
- Ali Al-Haj. Combined dwt-dct digital image watermarking. *Journal of computer science*, 3(9): 740–746, 2007.
- Sina Alemohammad, Josue Casco-Rodriguez, Lorenzo Luzi, Ahmed Imtiaz Humayun, Hossein Babaei, Daniel LeJeune, Ali Siahkoobi, and Richard Baraniuk. Self-consuming generative models go MAD. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=ShjMHfP50>.
- James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions. *Computer Science*. <https://cdn.openai.com/papers/dall-e-3.pdf>, 2(3):8, 2023.
- Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video generation models as world simulators. 2024. URL <https://openai.com/research/video-generation-models-as-world-simulators>.
- Miles Brundage, Shahar Avin, Jack Clark, Helen Toner, Peter Eckersley, Ben Garfinkel, Allan Dafoe, Paul Scharre, Thomas Zeitzoff, Bobby Filar, et al. The malicious use of artificial intelligence: Forecasting, prevention, and mitigation. *arXiv preprint arXiv:1802.07228*, 2018.
- Chin-Chen Chang, Piyu Tsai, and Chia-Chen Lin. Svd-based digital image watermarking scheme. *Pattern Recognition Letters*, 26(10):1577–1586, 2005.
- Siyi Chen, Zhang Huijie, Minzhe Guo, Yifu Lu, Peng Wang, and Qing Qu. Exploring low-dimensional subspaces in diffusion models for controllable image editing. In *Thirty-eighth Annual Conference on Neural Information Processing Systems (NeurIPS2024)*, 2024.
- Hai Ci, Pei Yang, Yiren Song, and Mike Zheng Shou. Ringid: Rethinking tree-ring watermarking for enhanced multi-key identification. *arXiv preprint arXiv:2404.14055*, 2024.
- Ingemar Cox, Matthew Miller, Jeffrey Bloom, Jessica Fridrich, and Ton Kalker. *Digital watermarking and steganography*. Morgan kaufmann, 2007.
- Elvis Dohmatob, Yunzhen Feng, Pu Yang, Francois Charton, and Julia Kempe. A tale of tails: Model collapse as a change of scaling laws. In *Forty-first International Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=KVvku47shW>.
- Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first International Conference on Machine Learning*, 2024.
- Zach Evans, Julian D Parker, CJ Carr, Zack Zukowski, Josiah Taylor, and Jordi Pons. Long-form music generation with latent diffusion. *arXiv preprint arXiv:2404.10301*, 2024.
- Pierre Fernandez, Guillaume Couairon, Hervé Jégou, Matthijs Douze, and Teddy Furon. The stable signature: Rooting watermarks in latent diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 22466–22477, 2023.
- Shi Fu, Sen Zhang, Yingjie Wang, Xinmei Tian, and Dacheng Tao. Towards theoretical understandings of self-consuming generative models. In *Forty-first International Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=aw6L8sB2Ts>.
- Elizabeth Gibney. Ai models fed ai-generated data quickly spew nonsense. *Nature*, 632(8023): 18–19, 2024.

-
- Josh A Goldstein, Girish Sastry, Micah Musser, Renee DiResta, Matthew Gentzel, and Katerina Sedova. Generative language models and automated influence operations: Emerging threats and potential mitigations. *arXiv preprint arXiv:2301.04246*, 2023.
- Xiangming Gu, Chao Du, Tianyu Pang, Chongxuan Li, Min Lin, and Ye Wang. On memorization in diffusion models. *arXiv preprint arXiv:2310.02664*, 2023.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.
- Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- Bernd Jähne. *Digital image processing*. Springer Science & Business Media, 2005.
- Hamidreza Kamkari, Brendan Leigh Ross, Rasa Hosseinzadeh, Jesse C Cresswell, and Gabriel Loaiza-Ganem. A geometric view of data complexity: Efficient local intrinsic dimension estimation with diffusion models. In *Thirty-eighth Annual Conference on Neural Information Processing Systems (NeurIPS2024)*, 2024.
- Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. *Advances in Neural Information Processing Systems*, 35:26565–26577, 2022.
- Jae-Eun Lee, Young-Ho Seo, and Dong-Wook Kim. Convolutional neural network-based digital image watermarking adaptive to the resolution of image and watermark. *Applied Sciences*, 10 (19):6854, 2020.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pp. 740–755. Springer, 2014.
- Junxiu Liu, Jiadong Huang, Yuling Luo, Lvchen Cao, Su Yang, Duqu Wei, and Ronglong Zhou. An optimized image watermarking method based on hd and svd in dwt domain. *IEEE Access*, 7: 80849–80860, 2019.
- Gabriel Loaiza-Ganem, Brendan Leigh Ross, Rasa Hosseinzadeh, Anthony L. Caterini, and Jesse C. Cresswell. Deep generative models through the lens of the manifold hypothesis: A survey and new connections. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL <https://openreview.net/forum?id=a90WpmSi0I>. Survey Certification, Expert Certification.
- Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *Advances in Neural Information Processing Systems*, 35:5775–5787, 2022.
- Calvin Luo. Understanding diffusion models: A unified perspective. *arXiv preprint arXiv:2208.11970*, 2022.
- Hila Manor and Tomer Michaeli. On the posterior distribution in denoising: Application to uncertainty quantification. In *The Twelfth International Conference on Learning Representations*, 2024a. URL <https://openreview.net/forum?id=adSGeugiuJ>.
- Hila Manor and Tomer Michaeli. Zero-shot unsupervised and text-based audio editing using DDPM inversion. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 34603–34629. PMLR, 21–27 Jul 2024b. URL <https://proceedings.mlr.press/v235/manor24a.html>.

-
- KA Navas, Mathews Cheriyan Ajay, M Lekshmi, Tampy S Archana, and M Sasikumar. Dwt-dct-svd based watermarking. In *2008 3rd international conference on communication systems software and middleware and workshops (COMSWARE'08)*, pp. 271–274. IEEE, 2008.
- Sandu Popescu, Anthony J Short, and Andreas Winter. Entanglement and the foundations of statistical mechanics. *Nature Physics*, 2(11):754–758, 2006.
- 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 *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- 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.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems*, 35:36479–36494, 2022.
- Ilya Shumailov, Zakhar Shumaylov, Yiren Zhao, Nicolas Papernot, Ross Anderson, and Yarin Gal. Ai models collapse when trained on recursively generated data. *Nature*, 631(8022):755–759, 2024.
- Vassilios Solachidis and Ioannis Pitas. Circularly symmetric watermark embedding in 2-d dft domain. *IEEE transactions on image processing*, 10(11):1741–1753, 2001.
- Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Diffusion art or digital forgery? investigating data replication in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6048–6058, 2023a.
- Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Understanding and mitigating copying in diffusion models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023b. URL <https://openreview.net/forum?id=HtMXRGbUMt>.
- Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2021a. URL <https://openreview.net/forum?id=StlgiaRCHLP>.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021b. URL <https://openreview.net/forum?id=PXTIG12RRHS>.
- Jan Pawel Stanczuk, Georgios Batzolis, Teo Deveney, and Carola-Bibiane Schönlieb. Diffusion models encode the intrinsic dimension of data manifolds. In *Forty-first International Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=a0XiA6v256>.
- Wei Ren Tan, Chee Seng Chan, Hernan Aguirre, and Kiyoshi Tanaka. Improved artgan for conditional synthesis of natural image and artwork. *IEEE Transactions on Image Processing*, 28(1): 394–409, 2019. doi: 10.1109/TIP.2018.2866698. URL <https://doi.org/10.1109/TIP.2018.2866698>.
- Peng Wang, Huijie Zhang, Zekai Zhang, Siyi Chen, Yi Ma, and Qing Qu. Diffusion models learn low-dimensional distributions via subspace clustering. *arXiv preprint arXiv:2409.02426*, 2024.

-
- Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004.
- Zijie J. Wang, Evan Montoya, David Munechika, Haoyang Yang, Benjamin Hoover, and Duen Horng Chau. Large-scale prompt gallery dataset for text-to-image generative models. *arXiv:2210.14896 [cs]*, 2022. URL <https://arxiv.org/abs/2210.14896>.
- Yuxin Wen, John Kirchenbauer, Jonas Geiping, and Tom Goldstein. Tree-rings watermarks: Invisible fingerprints for diffusion images. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023a. URL <https://openreview.net/forum?id=Z57JrmubNL>.
- Yuxin Wen, Yuchen Liu, Chen Chen, and Lingjuan Lyu. Detecting, explaining, and mitigating memorization in diffusion models. In *The Twelfth International Conference on Learning Representations*, 2023b.
- Zijin Yang, Kai Zeng, Kejiang Chen, Han Fang, Weiming Zhang, and Nenghai Yu. Gaussian shading: Provable performance-lossless image watermarking for diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12162–12171, 2024.
- Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and Yejin Choi. Defending against neural fake news. *Advances in neural information processing systems*, 32, 2019.
- Benjamin J Zhang, Siting Liu, Wuchen Li, Markos A Katsoulakis, and Stanley J Osher. Wasserstein proximal operators describe score-based generative models and resolve memorization. *arXiv preprint arXiv:2402.06162*, 2024a.
- Huijie Zhang, Jinfan Zhou, Yifu Lu, Minzhe Guo, Peng Wang, Liye Shen, and Qing Qu. The emergence of reproducibility and consistency in diffusion models. In *Forty-first International Conference on Machine Learning*, 2024b. URL <https://openreview.net/forum?id=HsliOqZkc0>.
- Kevin Alex Zhang, Lei Xu, Alfredo Cuesta-Infante, and Kalyan Veeramachaneni. Robust invisible video watermarking with attention. *arXiv preprint arXiv:1909.01285*, 2019.
- Lijun Zhang, Xiao Liu, Antoni Viro Martin, Cindy Xiong Bearfield, Yuriy Brun, and Hui Guan. Robust image watermarking using stable diffusion, 2024c. URL <https://arxiv.org/abs/2401.04247>.
- Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 586–595, 2018.
- Jiren Zhu, Russell Kaplan, Justin Johnson, and Li Fei-Fei. Hidden: Hiding data with deep networks. In *European Conference on Computer Vision*, 2018.

A RELATED WORK

A.1 IMAGE WATERMARKING

Image watermarking has long been a crucial method for protecting intellectual property in computer vision (Cox et al., 2007; Solachidis & Pitas, 2001; Chang et al., 2005; Liu et al., 2019). Traditional techniques primarily focus on user-side watermarking, where watermarks are embedded into images post-generation. These methods (Al-Haj, 2007; Navas et al., 2008) typically operate in the frequency domain to ensure the watermarks are imperceptible. However, such watermarks remain vulnerable to adversarial attacks and can become undetectable after applying simple image manipulations like blurring.

Early deep learning-based approaches to watermarking (Zhang et al., 2024c; Fernandez et al., 2023; Ahmadi et al., 2020; Lee et al., 2020; Zhu et al., 2018) leveraged neural networks to embed watermarks. While these methods improved robustness and imperceptibility, they often suffer from high

computational costs during fine-tuning and lack flexibility. Each new watermark requires additional fine-tuning or retraining, limiting their practicality.

More recently, diffusion model-based watermarking techniques have gained attraction due to their ability to seamlessly integrate watermarks during the generative process without incurring extra computational costs. Techniques such as Wen et al. (2023a); Yang et al. (2024); Ci et al. (2024) embed watermarks directly into the initial noise and retrieve the watermark by reversing the diffusion process. These methods enhance robustness and invisibility but are typically restricted to server-side watermarking, requiring access to the initial random seed. Moreover, the watermarks introduced by Wen et al. (2023a); Ci et al. (2024) significantly alter the data distribution, leading to variance towards watermarks in generated outputs (as shown in Figure 1).

In contrast to Wen et al. (2023a); Ci et al. (2024), our proposed shallow diffuse disentangles the watermark embedding from the generation process by leveraging the high-dimensional null space. This approach, both empirically and theoretically validated, significantly improves watermark consistency and robustness. To the best of our knowledge, this is the first method that supports watermark embedding for both server-side and user-side applications while maintaining high robustness and consistency.

A.2 LOW-DIMENSIONAL SUBSPACE IN DIFFUSION MODEL

In recent years, there has been growing interest in understanding deep generative models through the lens of the manifold hypothesis (Loaiza-Ganem et al., 2024). This hypothesis suggests that high-dimensional real-world data actually lies in latent manifolds with a low intrinsic dimension. Focusing on diffusion models, Stanczuk et al. (2024) empirically and theoretically shows that the approximated score function (the gradient of the log density of a noise-corrupted data distribution) in diffusion models is orthogonal to a low-dimensional subspace. Building on this, Wang et al. (2024); Chen et al. (2024) find that the estimated posterior mean from diffusion models lies within this low-dimensional space. Additionally, Chen et al. (2024) discovers strong local linearity within the space, suggesting that it can be locally approximated by a linear subspace. This observation motivates our Assumption 1, where we assume the estimated posterior mean lies in a low-dimensional subspace.

Building upon these findings, Stanczuk et al. (2024); Kamkari et al. (2024) introduce a local intrinsic dimension estimator, while Loaiza-Ganem et al. (2024) proposes a method for detecting out-of-domain data. Wang et al. (2024) offers theoretical insights into how diffusion model training transitions from memorization to generalization, and Chen et al. (2024); Manor & Michaeli (2024b) explores the semantic basis of the subspace to achieve disentangled image editing. Unlike these previous works, our approach leverages the low-dimensional subspace for watermarking, where both empirical and theoretical evidence demonstrates that this subspace enhances robustness and consistency.

B CHANNEL AVERAGING

B.1 TECHNIQUE DETAILS

Natural images have multiple channels denoted by C . Instead of applying watermark $\lambda\Delta$ to all channels of \mathbf{x}_t , we can apply the watermark to a specific channel c to make it even more invisible and robust. For this consideration, let us reshape the image \mathbf{x}_t and the watermark $\Delta\mathbf{x}$ into the form $\mathbf{x}_t \in \mathbb{R}^{H \times W \times C}$, $\lambda\Delta\mathbf{x} \in \mathbb{R}^{H \times W \times C}$, where H , W , and C represent the height, width, and channel dimensions for the image, respectively. These dimensions satisfy $HWC = d$.

Denote $[\mathbf{x}_t]_i \in \mathbb{R}^{H \times W}$ as the i th channel of \mathbf{x}_t , with $i \in [C]$. Thus $[\mathbf{x}_t^{\mathcal{W}}]_c = [\mathbf{x}_t]_c + [\lambda\Delta\mathbf{x}]_c$ and $[\mathbf{x}_t^{\mathcal{W}}]_i = [\mathbf{x}_t]_i$ for $i \neq c$. For the watermark in Equation (3), the channel averaging is defined

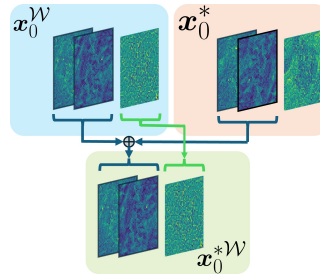


Figure 5: Illustration of channel average

as:

$$[\mathbf{x}_0^{*\mathcal{W}}]_i = \text{ChannelAverage}(\mathbf{x}_0^{\mathcal{W}}, \mathbf{x}_0^*), \quad (9)$$

$$= \begin{cases} [\mathbf{x}_0^{\mathcal{W}}]_i, i = c \\ (1 - \gamma)[\mathbf{x}_0^{\mathcal{W}}]_i + \gamma[\mathbf{x}_0^*]_i, i \neq c \end{cases}, \quad (10)$$

where we applied $\gamma = 1$. In our experiments, we found that we can increase both imperceptibility and robustness by further employing this simple approach. See our ablation study in Appendix C.5 for a more detailed analysis.

C ADDITIONAL EXPERIMENTS

C.1 DETAILS ABOUT EXPERIMENT SETTINGS

Baseline For the server scenario, we select the following methods as baselines: DWtDct Cox et al. (2007), DwtDctSvd Cox et al. (2007), RivaGAN Zhang et al. (2019), Stable Signature Fernandez et al. (2023), Tree-Ring Watermarks Wen et al. (2023a), RingId Ci et al. (2024), and Gaussian Shading Yang et al. (2024). In the user scenario, we adopt the same baseline methods, except for Stable Signature and Gaussian Shading, as these methods are not suitable for this setting.

Datasets We use Stable Diffusion 2.1 (Rombach et al., 2022) as the underlying model for our experiments, applying Shallow diffusion within its latent space. For the server scenario (Section 4.1), all diffusion-based methods are based on the same Stable Diffusion, with the original images \mathbf{x}_0 generated from identical initial seeds \mathbf{x}_T . Non-diffusion methods are applied to these same original images \mathbf{x}_0 in a post-watermarking process. A total of 5000 original images are generated for evaluation in this scenario. For the user scenario (Section 4.2), we utilize the MS-COCO Lin et al. (2014), WikiArt Tan et al. (2019), and DiffusionDB datasets Wang et al. (2022). The first two are real-world datasets, while DiffusionDB is a collection of diffusion model-generated images. From each dataset, we select 500 images for evaluation. For the remaining experiments in Appendix C.3, Section 4.3, Appendix C, we use the server scenario and sample 100 images for evaluation.

Metric To evaluate image consistency under the user scenario, we use peak signal-to-noise ratio (PSNR) Jähne (2005), structural similarity index measure (SSIM) Wang et al. (2004), and Learned Perceptual Image Patch Similarity (LPIPS) Zhang et al. (2018), comparing watermarked images to their original counterparts. In the server scenario, we assess the generation quality of the watermarked images using Contrastive Language-Image Pretraining Score (CLIP-Score) Radford et al. (2021) and Fréchet Inception Distance (FID) Heusel et al. (2017). To evaluate robustness, we vary the threshold η_0 and plot the true positive rate (TPR) against the false positive rate (FPR) for the receiver operating characteristic (ROC) curve. We use the area under the curve (AUC) and TPR when FPR = 0.01 (TPR @1% FPR) as robustness metrics. Robustness is evaluated both under clean conditions (no attacks) and with various attacks, including JPEG compression, Gaussian blurring (G.Blur), Gaussian noise, and color jitter. Details of these attacks are provided in Appendix C.2.

C.2 DETAILS ABOUT ATTACKS

In this work, we intensively tested our method on four different watermarking attacks, both in the server scenario and in the user scenario. These watermarking attacks represent the most common image distortion methods in real life, including

- JPEG compression with a compression rate of 25%
- Gaussian blurring (G.Blur) with an 8×8 filter size
- Gaussian noise (G.Noise) with $\sigma = 0.1$
- Color jitter with brightness factor uniformly ranges between 0 and 6

C.3 TRADE-OFF BETWEEN CONSISTENCY AND ROBUSTNESS

Figure 6 illustrates the trade-off between consistency and robustness for Shallow Diffuse and other baselines. As the radius of M increases, the watermark intensity λ also increases, reducing image

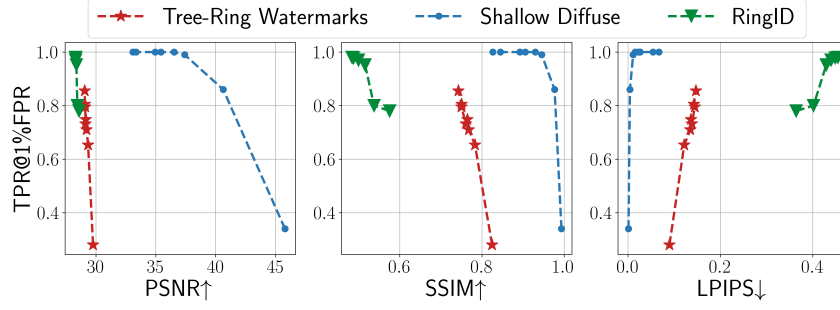


Figure 6: Trade-off between consistency and robustness for Tree-Ring Watermarks, RingID, and Shallow Diffuse.

Table 3: Ablation study of different watermark patterns.

Method & Dataset			PSNR ↑	SSIM ↑	LPIPS ↓	Average Watermarking Robustness (AUC ↑/TPR@1%FPR↑)
Frequency Region	Shape	Distribution				
Low	Circle	Zero	29.10	0.90	0.06	0.93/0.65
Low	Circle	Rand	29.37	0.92	0.05	0.92/0.25
Low	Circle	Rotational Rand	29.13	0.90	0.06	1.00/1.00
Low	Ring	Zero	36.20	0.95	0.02	0.78/0.35
Low	Ring	Rand	38.23	0.97	0.01	0.87/0.49
Low	Ring	Rotational Rand	35.23	0.93	0.02	0.99/0.98
High	Circle	Zero	38.3	0.96	0.01	0.80/0.34
High	Circle	Rand	42.3	0.98	0.004	0.86/0.35
High	Circle	Rotational Rand	38.0	0.94	0.01	1.00/1.00

consistency but improving robustness. By adjusting the radius of M , we plot the trade-off using PSNR, SSIM, and LPIPS against TPR@1%FPR. From Figure 6, curve of Shallow Diffuse is consistently above the curve of Tree-Ring Watermarks and RingID, demonstrating Shallow Diffuse’s better consistency at the same level of robustness.

C.4 ABLATION STUDY OF DIFFERENT WATERMARK PATTERNS

C.5 ABLATION STUDY OF CHANNEL AVERAGE

C.6 ABLATION STUDY OF WATERMARKING EMBEDDED CHANNEL.

D THEORETICAL JUSTIFICATION

In this section, we provide theoretical justifications for the consistency and the detectability of Shallow Diffuse introduced in Section 3 for unconditional diffusion models. First, we make the following assumptions on the watermark and the diffusion model process.

Assumption 1. Suppose the following hold for the PMP $\mathbf{f}_{\theta,t}(\mathbf{x}_t)$:

- **Linearity:** For any small t and $\Delta\mathbf{x} \in \mathbb{S}^{d-1}$, we always have

$$\mathbf{f}_{\theta,t}(\mathbf{x}_t + \lambda\Delta\mathbf{x}) = \mathbf{f}_{\theta,t}(\mathbf{x}_t) + \lambda\mathbf{J}_{\theta,t}(\mathbf{x}_t)\Delta\mathbf{x}.$$

- **L -Lipschitz continuous:** we assume that $\mathbf{f}_{\theta,t}(\mathbf{x})$ is a L -Lipschitz continuous at every t :

$$\|\mathbf{J}_{\theta,t}(\mathbf{x})\|_2 \leq L, \quad \forall \mathbf{x} \in \mathbb{R}^d.$$

It should be noted that our assumptions are mild. The L -Lipschitz continuity is a common assumption for analysis. The approximated linearity have been shown in (Chen et al., 2024) with the assumption of data distribution to be a mixture of low-rank Gaussians. Here, we assume the linearity to be exact for the ease of analysis, and it can be generalized to approximate linear case.

Now consider injecting a watermark $\lambda\Delta\mathbf{x}$ in Equation (3), where $\lambda > 0$ is a scaling factor and $\Delta\mathbf{x}$ is a random vector uniformly distributed on the unit hypersphere \mathbb{S}^{d-1} , i.e., $\Delta\mathbf{x} \sim \mathcal{U}(\mathbb{S}^{d-1})$. Then the following hold for the PMP $\mathbf{f}_{\theta,t}(\mathbf{x}_t)$.

Table 4: **Ablation study of channel average.**

Channel average intensity γ	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Watermarking Robustness (TPR@1%FPR \uparrow)				
				Clean	JPEG	G.Blur	G.Noise	Color Jitter
0	37.1103	0.941	0.0154	1.0000	1.0000	0.9971	1.0000	0.9584
1.0	36.6352	0.931	0.0151	1.0000	1.0000	1.0000	1.0000	1.0000

Table 5: **Ablation study of watermarking embedded channel.**

Watermark embedding channel	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Watermarking Robustness (TPR@1%FPR \uparrow)				
				Clean	JPEG	G.Blur	G.Noise	Color Jitter
0th	36.46	0.93	0.02	1.00	1.00	1.00	1.00	0.99
1th	36.57	0.93	0.02	1.00	1.00	1.00	1.00	0.99
2th	36.13	0.92	0.02	1.00	1.00	1.00	1.00	1.00
3th	36.64	0.93	0.02	1.00	1.00	1.00	1.00	1.00
1th + 2th + 3th	33.19	0.83	0.05	1.00	1.00	1.00	1.00	0.95

Theorem 1 (Consistency of the watermarks). *Suppose Assumption 1 holds and $\Delta\mathbf{x} \sim \mathcal{U}(\mathbb{S}^{d-1})$. Let us define $\hat{\mathbf{x}}_{0,t}^{\mathcal{W}} := \mathbf{f}_{\theta,t}(\mathbf{x}_t + \lambda\Delta\mathbf{x})$, $\hat{\mathbf{x}}_{0,t} := \mathbf{f}_{\theta,t}(\mathbf{x}_t)$. The ℓ_2 -norm distance between $\hat{\mathbf{x}}_{0,t}^{\mathcal{W}}$ and $\hat{\mathbf{x}}_{0,t}$ can be bounded by:*

$$\|\hat{\mathbf{x}}_{0,t}^{\mathcal{W}} - \hat{\mathbf{x}}_{0,t}\|_2 \leq \lambda L h(r_t), \quad (11)$$

with probability at least $1 - r_t^{-1}$. Here, $h(r_t) = \sqrt{\frac{r_t}{d}} + \sqrt{\frac{18\pi^3}{d-2} \log(2r_t)}$.

Our Theorem 1 guarantees that adding the watermark $\lambda\Delta\mathbf{x}$ would only change the estimation by an amount of $\lambda L h(r_t)$ with a constant probability. In particular, when r_t is small, it implies that the change in the prediction would be small. Given the relationship between PMP and DDIM in equation 1, the consistency also applies to the practical use. On the other hand, in the following we show that the injected watermark can be detected based upon the second term in Equation (6).

Theorem 2 (Detectability of the watermarks). *Suppose Assumption 1 holds and $\Delta\mathbf{x} \sim \mathcal{U}(\mathbb{S}^{d-1})$. With $\mathbf{x}_t^{\mathcal{W}}$ given in Equation (3), define $\mathbf{x}_{t-1}^{\mathcal{W}} = \text{DDIM}(\mathbf{x}_t^{\mathcal{W}}, t-1)$ and $\tilde{\mathbf{x}}_t^{\mathcal{W}} = \text{DDIM-Inv}(\mathbf{x}_{t-1}^{\mathcal{W}}, t)$. The ℓ_2 -norm distance between $\tilde{\mathbf{x}}_t^{\mathcal{W}}$ and $\mathbf{x}_t^{\mathcal{W}}$ can be bounded by:*

$$\|\tilde{\mathbf{x}}_t^{\mathcal{W}} - \mathbf{x}_t^{\mathcal{W}}\|_2 \leq \lambda L (-g(\alpha_t, \alpha_{t-1}) + g(\alpha_{t-1}, \alpha_t) (1 - Lg(\alpha_t, \alpha_{t-1}))) h(\max\{r_{t-1}, r_t\}) \quad (12)$$

with probability at least $1 - r_t^{-1} - r_{t-1}^{-1}$. Here, $g(x, y) := \frac{\sqrt{1-y}\sqrt{x} - \sqrt{1-x}\sqrt{y}}{\sqrt{1-x}}$, $\forall x, y \in (0, 1)$.

Here $-g(\alpha_t, \alpha_{t-1}) + g(\alpha_{t-1}, \alpha_t) (1 - Lg(\alpha_t, \alpha_{t-1}))$ is a small number under the α_t designed for variance preserving (VP) noise scheduler Ho et al. (2020) and $h(\max\{r_{t-1}, r_t\})$ is small when r_t is small. This indicates that the difference between $\tilde{\mathbf{x}}_t^{\mathcal{W}}$ and $\mathbf{x}_t^{\mathcal{W}}$ is small when r_t is small and $\mathbf{x}_t^{\mathcal{W}}$ could be recovered by $\tilde{\mathbf{x}}_t^{\mathcal{W}}$ from one-step DDIM. Therefore, Theorem 2 implies that the injected watermark can be detected with constant probability.

E PROOFS IN SECTION D

E.1 PROOFS OF THEOREM 1

Proof of Theorem 1. According to Assumption 1, we have $\|\hat{\mathbf{x}}_{0,t}^{\mathcal{W}} - \hat{\mathbf{x}}_{0,t}\|_2^2 = \lambda \|\mathbf{J}_{\theta,t}(\mathbf{x}_t) \cdot \Delta\mathbf{x}\|_2^2$. From Levy’s Lemma proposed in Popescu et al. (2006), given function $\|\mathbf{J}_{\theta,t}(\mathbf{x}_t) \cdot \Delta\mathbf{x}\|_2^2 : \mathbb{S}^{d-1} \rightarrow \mathbb{R}$ we have:

$$\mathbb{P}(\left| \|\mathbf{J}_{\theta,t}(\mathbf{x}_t) \cdot \Delta\mathbf{x}\|_2^2 - \mathbb{E}[\|\mathbf{J}_{\theta,t}(\mathbf{x}_t) \cdot \Delta\mathbf{x}\|_2^2] \right| \geq \epsilon) \leq 2 \exp\left(\frac{-C(d-2)\epsilon^2}{L^2}\right),$$

given L to be the Lipschitz constant of $\|\mathbf{J}_{\theta,t}(\mathbf{x}_t)\|_2^2$ and C is a positive constant (which can be taken to be $C = (18\pi^3)^{-1}$). From Lemma 2 and Lemma 3, we have:

$$\mathbb{P}\left(\left|\|\mathbf{J}_{\theta,t}(\mathbf{x}_t) \cdot \Delta \mathbf{x}\|_2^2 - \frac{\|\mathbf{J}_{\theta,t}(\mathbf{x}_t)\|_F^2}{d}\right| \geq \epsilon\right) \leq 2 \exp\left(\frac{-(18\pi^3)^{-1}(d-2)\epsilon^2}{\|\mathbf{J}_{\theta,t}(\mathbf{x}_t)\|_2^4}\right).$$

Define $\frac{1}{r_t}$ as the desired probability level, set

$$\frac{1}{r_t} = 2 \exp\left(\frac{-(18\pi^3)^{-1}(d-2)\epsilon^2}{\|\mathbf{J}_{\theta,t}(\mathbf{x}_t)\|_2^4}\right),$$

Solving for ϵ :

$$\epsilon = \|\mathbf{J}_{\theta,t}(\mathbf{x}_t)\|_2^2 \sqrt{\frac{18\pi^3}{d-2} \log(2r_t)}.$$

Therefore, with probability $1 - \frac{1}{r_t}$, we have:

$$\begin{aligned} \|\hat{\mathbf{x}}_{0,t}^{\mathcal{W}} - \hat{\mathbf{x}}_{0,t}\|_2^2 &= \lambda^2 \|\mathbf{J}_{\theta,t}(\mathbf{x}_t) \cdot \Delta \mathbf{x}\|_2^2, \\ &\leq \frac{\lambda^2 \|\mathbf{J}_{\theta,t}(\mathbf{x}_t)\|_F^2}{d} + \lambda^2 \|\mathbf{J}_{\theta,t}(\mathbf{x}_t)\|_2^2 \sqrt{\frac{18\pi^3}{d-2} \log(2r_t)}, \\ &\leq \lambda^2 \|\mathbf{J}_{\theta,t}(\mathbf{x}_t)\|_2^2 \left(\frac{r_t}{d} + \sqrt{\frac{18\pi^3}{d-2} \log(2r_t)}\right), \\ &= \lambda^2 L^2 \left(\frac{r_t}{d} + \sqrt{\frac{18\pi^3}{d-2} \log(2r_t)}\right), \end{aligned}$$

where the last inequality is obtained from $\|\mathbf{J}_{\theta,t}(\mathbf{x}_t)\|_F^2 \leq r_t \|\mathbf{J}_{\theta,t}(\mathbf{x}_t)\|_2^2$. Therefore, with probability $1 - \frac{1}{r_t}$,

$$\|\hat{\mathbf{x}}_{0,t}^{\mathcal{W}} - \hat{\mathbf{x}}_{0,t}\|_2 \leq \lambda L \sqrt{\frac{r_t}{d} + \sqrt{\frac{18\pi^3}{d-2} \log(2r_t)}} = \lambda L h(r_t).$$

□

Proof of Theorem 2. According to Equation (1), one step of DDIM sampling at timestep t could be represented by PMP $\mathbf{f}_{\theta,t}(\mathbf{x}_t)$ as:

$$\mathbf{x}_{t-1} = \sqrt{\alpha_{t-1}} \mathbf{f}_{\theta,t}(\mathbf{x}_t) + \sqrt{1 - \alpha_{t-1}} \left(\frac{\mathbf{x}_t - \sqrt{\alpha_t} \mathbf{f}_{\theta,t}(\mathbf{x}_t)}{\sqrt{1 - \alpha_t}} \right), \quad (13)$$

$$= \sqrt{\frac{1 - \alpha_{t-1}}{1 - \alpha_t}} \mathbf{x}_t + \frac{\sqrt{1 - \alpha_t} \sqrt{\alpha_{t-1}} - \sqrt{1 - \alpha_{t-1}} \sqrt{\alpha_t}}{\sqrt{1 - \alpha_t}} \mathbf{f}_{\theta,t}(\mathbf{x}_t), \quad (14)$$

If we inject a watermark $\lambda \Delta \mathbf{x}$ to \mathbf{x}_t , so $\mathbf{x}_t^{\mathcal{W}} = \mathbf{x}_t + \lambda \Delta \mathbf{x}$. To solve $\mathbf{x}_{t-1}^{\mathcal{W}}$, we could plugging Equation (2) to Equation (14), we could obtain:

$$\mathbf{x}_{t-1}^{\mathcal{W}} = \sqrt{\frac{1 - \alpha_{t-1}}{1 - \alpha_t}} \mathbf{x}_t^{\mathcal{W}} + \frac{\sqrt{1 - \alpha_t} \sqrt{\alpha_{t-1}} - \sqrt{1 - \alpha_{t-1}} \sqrt{\alpha_t}}{\sqrt{1 - \alpha_t}} \mathbf{f}_{\theta,t}(\mathbf{x}_t^{\mathcal{W}}), \quad (15)$$

$$= \mathbf{x}_{t-1} + \sqrt{\frac{1 - \alpha_{t-1}}{1 - \alpha_t}} \lambda \Delta \mathbf{x} + \frac{\sqrt{1 - \alpha_t} \sqrt{\alpha_{t-1}} - \sqrt{1 - \alpha_{t-1}} \sqrt{\alpha_t}}{\sqrt{1 - \alpha_t}} \mathbf{J}_{\theta,t}(\mathbf{x}_t) \Delta \mathbf{x} \quad (16)$$

$$= \mathbf{x}_{t-1} + \lambda \underbrace{\left(\sqrt{\frac{1 - \alpha_{t-1}}{1 - \alpha_t}} \mathbf{I} + \frac{\sqrt{1 - \alpha_t} \sqrt{\alpha_{t-1}} - \sqrt{1 - \alpha_{t-1}} \sqrt{\alpha_t}}{\sqrt{1 - \alpha_t}} \mathbf{J}_{\theta,t}(\mathbf{x}_t) \right)}_{:= \mathbf{W}_t} \Delta \mathbf{x}, \quad (17)$$

One step DDIM Inverse sampling at timestep $t - 1$ could be represented by PMP $\mathbf{f}_{\theta,t}(\mathbf{x}_t)$ as:

$$\mathbf{x}_t = \sqrt{\frac{1 - \alpha_t}{1 - \alpha_{t-1}}} \mathbf{x}_{t-1} + \frac{\sqrt{1 - \alpha_{t-1}} \sqrt{\alpha_t} - \sqrt{1 - \alpha_t} \sqrt{\alpha_{t-1}}}{\sqrt{1 - \alpha_{t-1}}} \mathbf{f}_{\theta,t-1}(\mathbf{x}_{t-1}), \quad (18)$$

To detect the watermark, we apply one step DDIM Inverse on $\mathbf{x}_{t-1}^{\mathcal{W}}$ at timestep $t - 1$ to obtain $\tilde{\mathbf{x}}_t^{\mathcal{W}}$:

$$\begin{aligned} \tilde{\mathbf{x}}_t^{\mathcal{W}} &= \sqrt{\frac{1 - \alpha_t}{1 - \alpha_{t-1}}} \mathbf{x}_{t-1}^{\mathcal{W}} + \frac{\sqrt{1 - \alpha_{t-1}} \sqrt{\alpha_t} - \sqrt{1 - \alpha_t} \sqrt{\alpha_{t-1}}}{\sqrt{1 - \alpha_{t-1}}} \mathbf{f}_{\theta,t-1}(\mathbf{x}_{t-1}^{\mathcal{W}}), \\ &= \mathbf{x}_t + \lambda \underbrace{\left(\sqrt{\frac{1 - \alpha_t}{1 - \alpha_{t-1}}} \mathbf{I} + \frac{\sqrt{1 - \alpha_{t-1}} \sqrt{\alpha_t} - \sqrt{1 - \alpha_t} \sqrt{\alpha_{t-1}}}{\sqrt{1 - \alpha_{t-1}}} \mathbf{J}_{\theta,t-1}(\mathbf{x}_{t-1}) \right)}_{:= \mathbf{W}_{t-1}} \mathbf{W}_t \Delta \mathbf{x}, \\ &= \mathbf{x}_t + \lambda \mathbf{W}_{t-1} \mathbf{W}_t \Delta \mathbf{x} = \mathbf{x}_t^{\mathcal{W}} + \lambda (\mathbf{W}_{t-1} \mathbf{W}_t - \mathbf{I}) \Delta \mathbf{x}. \end{aligned}$$

Therefore:

$$\begin{aligned} \|\tilde{\mathbf{x}}_t^{\mathcal{W}} - \mathbf{x}_t^{\mathcal{W}}\|_2 &= \lambda \|(\mathbf{W}_{t-1} \mathbf{W}_t - \mathbf{I}) \Delta \mathbf{x}\|_2, \\ &= \lambda \left\| \frac{\sqrt{1 - \alpha_{t-1}} \sqrt{\alpha_t} - \sqrt{1 - \alpha_t} \sqrt{\alpha_{t-1}}}{\sqrt{1 - \alpha_t}} \mathbf{J}_{\theta,t-1}(\mathbf{x}_{t-1}) \Delta \mathbf{x}, \right. \\ &\quad \left. + \frac{\sqrt{1 - \alpha_t} \sqrt{\alpha_{t-1}} - \sqrt{1 - \alpha_{t-1}} \sqrt{\alpha_t}}{\sqrt{1 - \alpha_{t-1}}} \mathbf{J}_{\theta,t}(\mathbf{x}_t) \Delta \mathbf{x}, \right. \\ &\quad \left. - \frac{(\sqrt{1 - \alpha_t} \sqrt{\alpha_{t-1}} - \sqrt{1 - \alpha_{t-1}} \sqrt{\alpha_t})^2}{\sqrt{1 - \alpha_{t-1}} \sqrt{1 - \alpha_t}} \mathbf{J}_{\theta,t-1}(\mathbf{x}_{t-1}) \mathbf{J}_{\theta,t}(\mathbf{x}_t) \Delta \mathbf{x} \right\|_2, \\ &\leq -\lambda g(\alpha_t, \alpha_{t-1}) \|\mathbf{J}_{\theta,t-1}(\mathbf{x}_{t-1}) \Delta \mathbf{x}\|_2 + \lambda g(\alpha_{t-1}, \alpha_t) \|\mathbf{J}_{\theta,t}(\mathbf{x}_t) \Delta \mathbf{x}\|_2 \\ &\quad - \lambda g(\alpha_{t-1}, \alpha_t) g(\alpha_t, \alpha_{t-1}) \|\mathbf{J}_{\theta,t-1}(\mathbf{x}_{t-1}) \mathbf{J}_{\theta,t}(\mathbf{x}_t) \Delta \mathbf{x}\|_2, \\ &\leq -\lambda g(\alpha_t, \alpha_{t-1}) \|\mathbf{J}_{\theta,t-1}(\mathbf{x}_{t-1}) \Delta \mathbf{x}\|_2 \\ &\quad + \lambda g(\alpha_{t-1}, \alpha_t) (1 - g(\alpha_t, \alpha_{t-1}) L) \|\mathbf{J}_{\theta,t}(\mathbf{x}_t) \Delta \mathbf{x}\|_2, \\ &= -g(\alpha_t, \alpha_{t-1}) \|\hat{\mathbf{x}}_{0,t-1}^{\mathcal{W}} - \hat{\mathbf{x}}_{0,t-1}\|_2 \\ &\quad + g(\alpha_{t-1}, \alpha_t) (1 - g(\alpha_t, \alpha_{t-1}) L) \|\hat{\mathbf{x}}_{0,t}^{\mathcal{W}} - \hat{\mathbf{x}}_{0,t}\|_2, \end{aligned}$$

The first inequality holds because $g(\alpha_{t-1}, \alpha_t) < 0$ and $g(\alpha_t, \alpha_{t-1}) > 0$. The second inequality holds because $\|\mathbf{J}_{\theta,t-1}(\mathbf{x}_{t-1}) \mathbf{J}_{\theta,t}(\mathbf{x}_t) \Delta \mathbf{x}\|_2 \leq \|\mathbf{J}_{\theta,t-1}(\mathbf{x}_{t-1})\|_2 \|\mathbf{J}_{\theta,t}(\mathbf{x}_t) \Delta \mathbf{x}\|_2 \leq L \|\mathbf{J}_{\theta,t}(\mathbf{x}_t) \Delta \mathbf{x}\|_2$. From Theorem 1, with probability $1 - \frac{1}{r_{t-1}}$,

$$\|\hat{\mathbf{x}}_{0,t-1}^{\mathcal{W}} - \hat{\mathbf{x}}_{0,t-1}\|_2 \leq \lambda L h(r_{t-1}),$$

with probability $1 - \frac{1}{r_t}$,

$$\|\hat{\mathbf{x}}_{0,t}^{\mathcal{W}} - \hat{\mathbf{x}}_{0,t}\|_2 \leq \lambda L h(r_t),$$

Thus, from the union of bound, with a probability at least $1 - \frac{1}{r_t} - \frac{1}{r_{t-1}}$,

$$\begin{aligned} \|\tilde{\mathbf{x}}_t^{\mathcal{W}} - \mathbf{x}_t^{\mathcal{W}}\|_2 &\leq -\lambda L g(\alpha_t, \alpha_{t-1}) h(r_{t-1}) + \lambda L g(\alpha_{t-1}, \alpha_t) (1 - g(\alpha_t, \alpha_{t-1}) L) h(r_t) \\ &\leq \lambda L (-g(\alpha_t, \alpha_{t-1}) + g(\alpha_{t-1}, \alpha_t) (1 - L g(\alpha_t, \alpha_{t-1}))) h(\max\{r_{t-1}, r_t\}) \end{aligned}$$

□

F AUXILIARY RESULTS

Lemma 1. Given a unit vector \mathbf{v}_i with and $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$, we have

$$\mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)}[(\mathbf{v}_i^T \boldsymbol{\epsilon})^2 / \|\boldsymbol{\epsilon}\|_2^2] = \frac{1}{d}.$$

Proof of Lemma 1. Because $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$,

$$\mathbf{v}_i^T \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{v}_i^T \mathbf{0}, \mathbf{v}_i^T \mathbf{I}_d \mathbf{v}_i) = \mathcal{N}(\mathbf{v}_i^T \mathbf{0}, \mathbf{v}_i^T \mathbf{I}_d \mathbf{v}_i) = \mathcal{N}(0, 1), \quad (19)$$

Assume a set of d unit vecotrs $\{v_1, v_2, \dots, \mathbf{v}_i, \dots, v_d\}$ are orthogonormal and are basis of \mathbb{R}^d , similarly, we could show that $\forall j \in [d], X_j := \mathbf{v}_j^T \boldsymbol{\epsilon} \sim \mathcal{N}(0, 1)$. Therefore, we could rewrite $(\mathbf{v}_i^T \boldsymbol{\epsilon})^2 / \|\boldsymbol{\epsilon}\|_2^2$ as:

$$(\mathbf{v}_i^T \boldsymbol{\epsilon})^2 / \|\boldsymbol{\epsilon}\|_2^2 = \frac{(\mathbf{v}_i^T \boldsymbol{\epsilon})^2}{\|\sum_{k=1}^d v_k v_k^T \boldsymbol{\epsilon}\|_2^2}, \quad (20)$$

$$= \frac{(\mathbf{v}_i^T \boldsymbol{\epsilon})^2}{\sum_{k=1}^d (\mathbf{v}_k^T \boldsymbol{\epsilon})^2}, \quad (21)$$

$$= \frac{X_i^2}{\sum_{k=1}^d X_k^2}. \quad (22)$$

Let $Y_i := \frac{X_i^2}{\sum_{j=1}^d X_j^2}$. Because $\forall j \in [d], X_j := \mathbf{v}_j^T \boldsymbol{\epsilon} \sim \mathcal{N}(0, 1), \forall j \in [d], Y_j$ has the same distribution. Additionally, $\sum_{j=1}^d Y_j = 1$. So:

$$\mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)}\left[\frac{(\mathbf{v}_i^T \boldsymbol{\epsilon})^2}{\|\boldsymbol{\epsilon}\|_2^2}\right] = \mathbb{E}[Y_i] = \frac{1}{d} \mathbb{E}\left[\sum_{j=1}^d Y_j\right] = \frac{1}{d}.$$

□

Lemma 2. Given a matrix $\mathbf{J} \in \mathbb{R}^{d \times d}$ with $\text{rank}(\mathbf{J}) = r$. Given \mathbf{x} which is uniformly sampled on the unit hypersphere \mathbb{S}^{d-1} , we have:

$$\mathbb{E}_{\mathbf{x}} [\|\mathbf{J}\mathbf{x}\|_2^2] = \frac{\|\mathbf{J}\|_F^2}{d}.$$

Proof of Lemma 2. Let's define the singular value decomposition of $\mathbf{J} = \mathbf{U}\Sigma\mathbf{V}^T$ with $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_r, 0, \dots, 0)$. Therefore, $\mathbb{E}_{\mathbf{x}} [\|\mathbf{J}\mathbf{x}\|_2^2] = \mathbb{E}_{\mathbf{x}} [\|\mathbf{U}\Sigma\mathbf{V}^T\mathbf{x}\|_2^2] = \mathbb{E}_{\mathbf{z}} [\|\Sigma\mathbf{z}\|_2^2]$ where $\mathbf{z} := \mathbf{V}^T\mathbf{x}$ is uniformly sampled on the unit hypersphere \mathbb{S}^{d-1} . Thus, we have:

$$\begin{aligned} \mathbb{E}_{\mathbf{z}} [\|\Sigma\mathbf{z}\|_2^2] &= \mathbb{E}_{\mathbf{z}} \left[\left\| \sum_{i=1}^r \sigma_i \mathbf{e}_i^T \mathbf{z} \right\|_2^2 \right], \\ &= \mathbb{E}_{\mathbf{z}} \left[\sum_{i=1}^r \sigma_i^2 \|\mathbf{e}_i^T \mathbf{z}\|_2^2 \right], \\ &= \sum_{i=1}^r \sigma_i^2 \mathbb{E}_{\mathbf{z}} [\|\mathbf{e}_i^T \mathbf{z}\|_2^2] = \frac{\|\mathbf{J}\|_F^2}{d}, \end{aligned}$$

where \mathbf{e}_i is the standard basis with i -th element equals to 1. The second equality is because of independence between $\mathbf{e}_i^T \mathbf{z}$ and $\mathbf{e}_j^T \mathbf{z}$. The fourth equality is from Lemma 1. □

Lemma 3. Given function $f(\mathbf{x}) = \|\mathbf{J}\mathbf{x}\|_2^2$, the lipschitz constant L_f of function $f(\mathbf{x})$ is:

$$L_f = 2\|\mathbf{J}\|_2^2.$$

Proof of Lemma 3. The jacobian of $f(\mathbf{x})$ is:

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = 2\mathbf{J}^T \mathbf{J} \mathbf{x},$$

Therefore, the lipschitz constant L follows:

$$L_f = \sup_{\mathbf{x} \in \mathbb{S}^{d-1}} \|\nabla_{\mathbf{x}} f(\mathbf{x})\|_2 = 2 \sup_{\mathbf{x} \in \mathbb{S}^{d-1}} \|\mathbf{J}^T \mathbf{J} \mathbf{x}\|_2 = \|\mathbf{J}^T \mathbf{J}\|_2 = \|\mathbf{J}\|_2^2$$

□