

# SHALLOW DIFFUSE: ROBUST AND INVISIBLE WATERMARKING THROUGH LOW-DIMENSIONAL SUBSPACES IN DIFFUSION MODELS

**Anonymous authors**

Paper under double-blind review

## 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 watermarking into the sampling process of diffusion models. For example, the work (Ci et al., 2024;

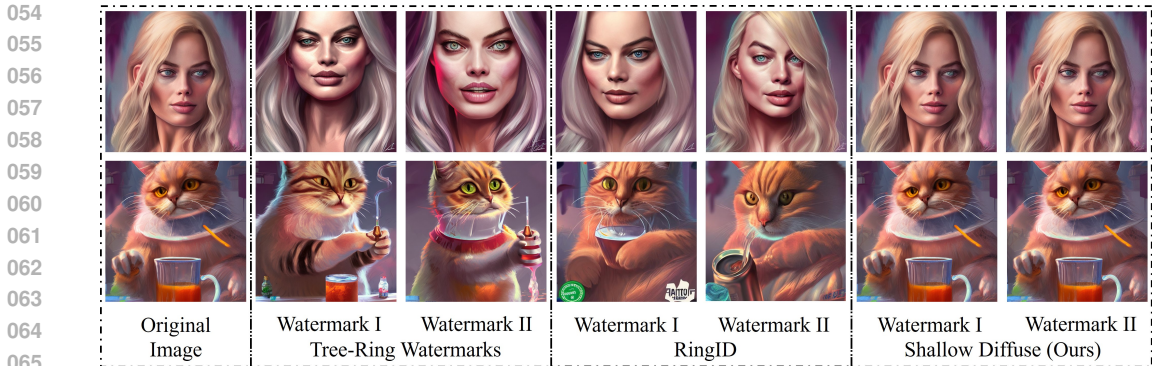


Figure 1: **Sampling variance of Tree-Ring Watermarks, RingID and Shallow Diffuse.** On the left are the original images, and on the right are the corresponding watermarked images generated using three different techniques: Tree-Ring (Wen et al., 2023a), RingID (Ci et al., 2024), and Shallow Diffuse.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

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 5a and Figure 5b.
- **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. (2022a); Karras et al. (2022) have proposed various methods to approximate the reverse process of diffusion models. Typically, these methods involve estimating the noise  $\epsilon_t$  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  $\epsilon_t$  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

---

```

162 1: Inject watermark:
163
164 2: Input: original image  $\mathbf{x}_0$  for the user scenario (initial random seed  $\mathbf{x}_T$  for the server scenario), watermark
165    $\lambda\Delta\mathbf{x}$ , embedding timestep  $t$ ,
166 3: Output: watermarked image  $\mathbf{x}_0^{*\mathcal{W}}$ ,
167 4: if user scenario then
168 5:    $\mathbf{x}_t = \text{DDIM-Inv}(\mathbf{x}_0, t)$ 
169 6: else server scenario
170 7:    $\mathbf{x}_t = \text{DDIM}(\mathbf{x}_T, t)$ 
171 8: end if
172 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
173 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
174 11: Return:  $\mathbf{x}_0^{*\mathcal{W}}$ 
175
176 12:
177 13: Detect watermark:
178 14: Input: Attacked image  $\bar{\mathbf{x}}_0^{\mathcal{W}}$ , watermark  $\lambda\Delta\mathbf{x}$ , embedding timestep  $t$ ,
179 15: Output: Distance score  $\eta$ ,
180 16:  $\bar{\mathbf{x}}_t^{\mathcal{W}} \leftarrow \text{DDIM-Inv}(\bar{\mathbf{x}}_0^{\mathcal{W}}, t)$ 
181 17:  $\eta = \text{Detector}(\bar{\mathbf{x}}_t^{\mathcal{W}}, \lambda\Delta\mathbf{x})$ 
182 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 5.4. 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:

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

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 Section 4, 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. **Specifically, the radius of the circle in mask  $\mathbf{M}$  is 8**, In the following, we discuss the design of  $\mathbf{M}$  and  $\mathbf{W}$  in detail.

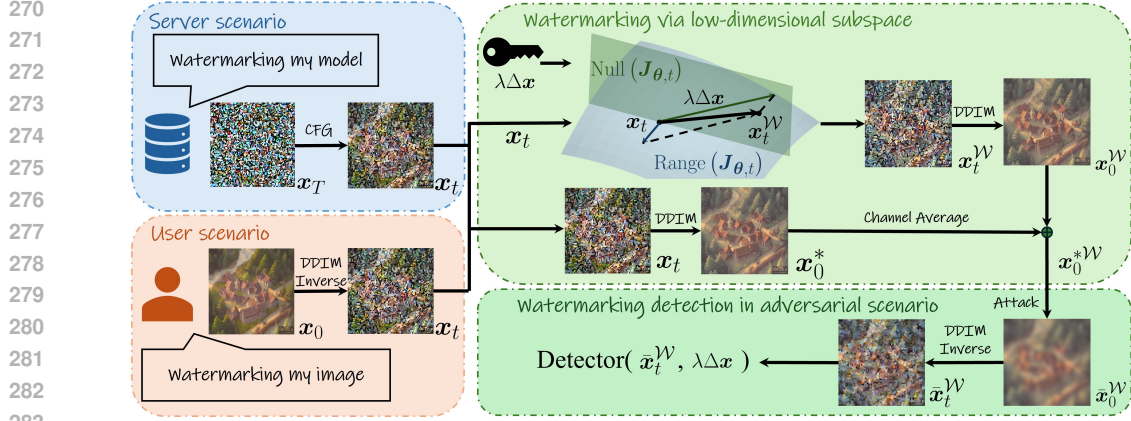


Figure 3: Overview of Shallow Diffuse for T2I diffusion models.

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 7.

**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  $\eta_0$ , non-watermarked images will have  $\eta > \eta_0$  and watermarked images will have  $\eta < \eta_0$ . Theoretically, the derivation of the p-value  $\eta$  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.

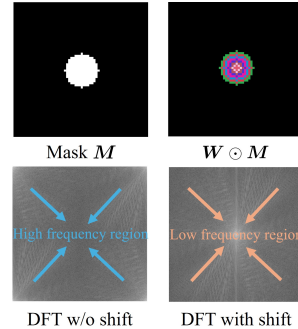


Figure 2: Illustration of watermark patterns.

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

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

## 4 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.,  $\Delta\mathbf{x} \sim \text{U}(\mathbb{S}^{d-1})$ . Then the following hold for the PMP  $\mathbf{f}_{\theta,t}(\mathbf{x}_t)$ .

**Theorem 1** (Consistency of the watermarks). *Suppose Assumption 1 holds and  $\Delta\mathbf{x} \sim \text{U}(\mathbb{S}^{d-1})$ . 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  $\ell_2$ -norm distance between  $\hat{\mathbf{x}}_{0,t}^{\mathcal{W}}$  and  $\hat{\mathbf{x}}_{0,t}$  can be bounded by:*

$$\|\hat{\mathbf{x}}_{0,t}^{\mathcal{W}} - \hat{\mathbf{x}}_{0,t}\|_2 \leq \lambda L h(r_t), \quad (9)$$

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  $\Delta\mathbf{x} \sim \text{U}(\mathbb{S}^{d-1})$ . 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  $\hat{\mathbf{x}}_t^{\mathcal{W}} = \text{DDIM-Inv}(\mathbf{x}_{t-1}^{\mathcal{W}}, t)$ . The  $\ell_2$ -norm distance between  $\hat{\mathbf{x}}_t^{\mathcal{W}}$  and  $\mathbf{x}_t^{\mathcal{W}}$  can be bounded by:*

$$\|\hat{\mathbf{x}}_t^{\mathcal{W}} - \mathbf{x}_t^{\mathcal{W}}\|_2 \leq \lambda L (-g(\alpha_t, \alpha_{t-1}) + g(\alpha_{t-1}, \alpha_t) (1 - Lg(\alpha_t, \alpha_{t-1}))) h(\max\{r_{t-1}, r_t\}) \quad (10)$$

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  $\alpha_t$  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  $\tilde{x}_t^{\mathcal{W}}$  and  $x_t^{\mathcal{W}}$  is small when  $r_t$  is small and  $x_t^{\mathcal{W}}$  could be recovered by  $\tilde{x}_t^{\mathcal{W}}$  from one-step DDIM. Therefore, Theorem 2 implies that the injected watermark can be detected with constant probability.

## 5 EXPERIMENTS

In this section, we present a comprehensive set of experiments to demonstrate the robustness and consistency of *Shallow-Diffuse* across various datasets. We begin by highlighting its performance in terms of robustness and consistency in both the server scenario (Section 5.1) and the user scenario (Section 5.2). Additionally, we compare Shallow Diffuse with other related works in the trade-off between robustness and consistency, as detailed in Section 5.3. Moreover, we investigate the effect of timestep  $t$  on both robustness and consistency, with results presented in Section 5.4. We further explore the multi-key identification experiments in Appendix C.2. Lastly, we provide an ablation study on watermark pattern design (Appendix C.3), channel averaging (Appendix C.4), watermarking embedded channel (Appendix C.5), and sampling method (Appendix C.6).

**Baseline** For the server scenario, we select the following non-diffusion-based method: DWtDet Cox et al. (2007), DwtDetSvd Cox et al. (2007), RivaGAN Zhang et al. (2019), **StegaStamp** Tancik et al. (2020); and diffusion-based method: 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-base** (Rombach et al., 2022) as the underlying model for our experiments, applying Shallow diffusion within its latent space. For the server scenario (Section 5.1), all diffusion-based methods are based on the same Stable Diffusion, with the original images  $x_0$  generated from identical initial seeds  $x_T$ . Non-diffusion methods are applied to these same original images  $x_0$  in a post-watermarking process. A total of 5000 original images are generated for evaluation in this scenario. For the user scenario (Section 5.2), 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 Section 5.3, Section 5.4, 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  $\eta_0$  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, Gaussian noise, and color jitter, **Resize and restore**, **Random drop**, **median blurring**, **diffusion purification** Nie et al. (2022), **VAE-based image compression models** Cheng et al. (2020); **Ballé et al. (2018)** and **stable diffusion-based image regeneration** Zhao et al. (2023b). We report the average robustness of these attacks in the main paper. Detailed settings and experiment results of these attacks are provided in Appendix C.1.

### 5.1 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



Table 1: Generation quality and watermark robustness under the server scenario.

Method	Generation Quality		Watermark Robustness (AUC $\uparrow$ /TPR@1%FPR $\uparrow$ )	
	CLIP-Score $\uparrow$	FID $\downarrow$	Clean	Adversarial Average
Stable Diffusion w/o WM	0.3286	25.56	-	-
DwtDct	0.3298	25.73	0.97/0.85	0.61/0.18
DwtDctSvd	0.3291	26.00	<b>1.00/1.00</b>	0.79/0.46
RivaGAN	0.3252	24.60	1.00/0.99	0.85/0.57
Stegastamp	0.3552	<b>24.59</b>	<b>1.00/1.00</b>	0.97/0.87
Stable Signature	0.3622	30.86	<b>1.00/1.00</b>	0.83/0.44
Tree-Ring Watermarks	0.3310	25.82	<b>1.00/1.00</b>	0.98/0.87
RingID	0.3285	27.13	<b>1.00/1.00</b>	<b>1.00/1.00</b>
Gaussian Shading	<b>0.3631</b>	26.17	<b>1.00/1.00</b>	<b>1.00/1.00</b>
<b>Shallow Diffuse (ours)</b>	0.3285	25.58	<b>1.00/1.00</b>	<b>1.00/1.00</b>

Table 2: Generation consistency and watermark robustness under the user scenario.

Dataset	Method	Generation Consistency			Watermark Robustness (AUC $\uparrow$ /TPR@1%FPR $\uparrow$ )	
		PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	Clean	Adversarial Avg.
COCO	Stable Diffusion w/o WM	32.28	0.78	0.06	-	-
	DwtDct	37.88	0.97	<b>0.02</b>	0.98/0.83	0.61/0.19
	DwtDctSvd	38.06	<b>0.98</b>	<b>0.02</b>	<b>1.00/1.00</b>	0.79/0.48
	RivaGAN	<b>40.57</b>	<b>0.98</b>	0.04	<b>1.00/1.00</b>	0.87/0.61
	Stegastamp	31.88	0.86	0.08	<b>1.00/1.00</b>	0.96/0.83
	Tree-Ring Watermarks	28.22	0.51	0.41	<b>1.00/1.00</b>	0.99/0.93
	RingID	28.22	0.38	0.61	<b>1.00/1.00</b>	<b>1.00/0.99</b>
	<b>Shallow Diffuse (ours)</b>	32.11	0.77	0.06	<b>1.00/1.00</b>	1.00/0.98
DiffusionDB	Stable Diffusion w/o WM	33.42	0.85	0.03	-	-
	DwtDct	37.77	0.96	<b>0.02</b>	0.96/0.76	0.61/0.18
	DwtDctSvd	37.84	0.97	<b>0.02</b>	<b>1.00/1.00</b>	0.79/0.46
	RivaGAN	<b>40.6</b>	<b>0.98</b>	0.04	1.00/0.98	0.85/0.57
	Stegastamp	32.03	0.85	0.08	<b>1.00/1.00</b>	0.96/0.84
	Tree-Ring Watermarks	28.3	0.62	0.29	<b>1.00/1.00</b>	0.97/0.85
	RingID	27.9	0.21	0.77	<b>1.00/1.00</b>	<b>1.00/0.99</b>
	<b>Shallow Diffuse (ours)</b>	33.07	0.84	0.04	<b>1.00/1.00</b>	0.99/0.97
WikiArt	Stable Diffusion w/o WM	31.6	0.7	0.09	-	-
	DwtDct	38.84	0.97	0.02	0.96/0.75	0.60/0.18
	DwtDctSvd	39.14	0.98	0.02	<b>1.00/1.00</b>	0.78/0.48
	RivaGAN	<b>40.44</b>	<b>0.98</b>	<b>0.05</b>	<b>1.00/1.00</b>	0.87/0.60
	Stegastamp	31.62	0.85	0.09	<b>1.00/1.00</b>	0.95/0.75
	Tree-Ring Watermarks	28.24	0.53	0.34	<b>1.00/1.00</b>	0.97/0.92
	RingID	27.90	0.19	0.78	<b>1.00/1.00</b>	0.99/0.98
	<b>Shallow Diffuse (ours)</b>	31.4	0.68	0.10	<b>1.00/1.00</b>	<b>1.00/0.99</b>

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.

## 5.2 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  $\hat{x}_0 = \text{DDIM}(\text{DDIM-Inv}(\mathbf{x}_0, t, \emptyset), 0, \emptyset)$ , and measure the data consistency between  $\hat{x}_0$  and  $\mathbf{x}_0$  (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 non-diffusion-based methods, which are not affected by DDIM inversion errors, better image consistency is achievable. **However, as visualized in Figure 8, Shallow Diffuse also demonstrates strong generation consistency.** As for the robustness, Shallow Diffuse is comparable to RingID and outperforms all other methods in all three datasets. While RivaGAN achieves the best image consistency and comparable watermark 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.

## 5.3 TRADE-OFF BETWEEN CONSISTENCY AND ROBUSTNESS

Figure 4 illustrates the trade-off between consistency and robustness for Shallow Diffuse and other baselines. As the radius of  $M$  increases, the watermark intensity  $\lambda$  also increases, reducing image

486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

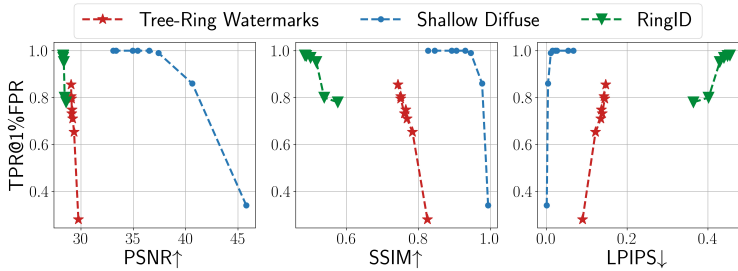


Figure 4: Trade-off between consistency and robustness for Tree-Ring Watermarks, RingID, and Shallow Diffuse.

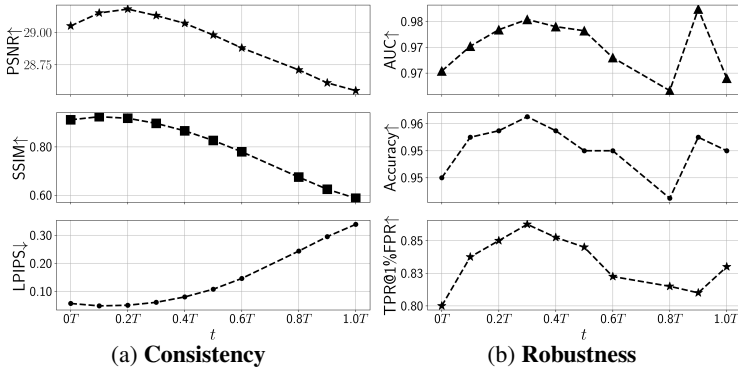


Figure 5: Ablation study of the watermark at different timestep  $t$ .

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 4, 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.

5.4 RELATION BETWEEN INJECTING TIMESTEP, CONSISTENCY AND ROBUSTNESS

Figure 5 shows the relationship between the watermark injection timestep  $t$  and both consistency and robustness<sup>1</sup>. 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 Section 4: 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  $J_{\theta,t}(x_t)$  remains low-rank but  $t$  is still below  $t^*$ . Another possible explanation is the gap between the image space and latent space in diffusion models. The rank curve in Chen et al. (2024) is evaluated for an image-space diffusion model, whereas Shallow Diffuse operates in the latent-space diffusion model (e.g., Stable Diffusion).

6 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.

<sup>1</sup>In 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.

## REFERENCES

- 540  
541  
542 Mahdi Ahmadi, Alireza Norouzi, Nader Karimi, Shadrokh Samavi, and Ali Emami. Redmark:  
543 Framework for residual diffusion watermarking based on deep networks. *Expert Systems with*  
544 *Applications*, 146:113157, 2020.
- 545 Ali Al-Haj. Combined dwt-dct digital image watermarking. *Journal of computer science*, 3(9):  
546 740–746, 2007.
- 547  
548 Sina Alemohammad, Josue Casco-Rodriguez, Lorenzo Luzi, Ahmed Imtiaz Humayun, Hossein  
549 Babaei, Daniel LeJeune, Ali Siahkoochi, and Richard Baraniuk. Self-consuming generative mod-  
550 els go MAD. In *The Twelfth International Conference on Learning Representations*, 2024. URL  
551 <https://openreview.net/forum?id=ShjMHfmPs0>.
- 552 Johannes Ballé, David Minnen, Saurabh Singh, Sung Jin Hwang, and Nick Johnston. Variational  
553 image compression with a scale hyperprior. In *International Conference on Learning Represen-*  
554 *tations*, 2018. URL <https://openreview.net/forum?id=rkcQFMZRb>.
- 555  
556 James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang  
557 Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions. *Computer*  
558 *Science*. <https://cdn.openai.com/papers/dall-e-3.pdf>, 2(3):8, 2023.
- 559 Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe  
560 Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video  
561 generation models as world simulators. 2024. URL [https://openai.com/research/](https://openai.com/research/video-generation-models-as-world-simulators)  
562 [video-generation-models-as-world-simulators](https://openai.com/research/video-generation-models-as-world-simulators).
- 563  
564 Miles Brundage, Shahar Avin, Jack Clark, Helen Toner, Peter Eckersley, Ben Garfinkel, Allan  
565 Dafoe, Paul Scharre, Thomas Zeitzoff, Bobby Filar, et al. The malicious use of artificial in-  
566 telligence: Forecasting, prevention, and mitigation. *arXiv preprint arXiv:1802.07228*, 2018.
- 567 Chin-Chen Chang, Piyu Tsai, and Chia-Chen Lin. Svd-based digital image watermarking scheme.  
568 *Pattern Recognition Letters*, 26(10):1577–1586, 2005.
- 569  
570 Siyi Chen, Zhang Huijie, Minzhe Guo, Yifu Lu, Peng Wang, and Qing Qu. Exploring low-  
571 dimensional subspaces in diffusion models for controllable image editing. In *Thirty-eighth Annual*  
572 *Conference on Neural Information Processing Systems (NeurIPS2024)*, 2024.
- 573 Zhengxue Cheng, Heming Sun, Masaru Takeuchi, and Jiro Katto. Learned image compression with  
574 discretized gaussian mixture likelihoods and attention modules. In *Proceedings of the IEEE/CVF*  
575 *conference on computer vision and pattern recognition*, pp. 7939–7948, 2020.
- 576  
577 Hai Ci, Pei Yang, Yiren Song, and Mike Zheng Shou. Ringid: Rethinking tree-ring watermarking  
578 for enhanced multi-key identification. *arXiv preprint arXiv:2404.14055*, 2024.
- 579 Ingemar Cox, Matthew Miller, Jeffrey Bloom, Jessica Fridrich, and Ton Kalker. *Digital watermark-*  
580 *ing and steganography*. Morgan kaufmann, 2007.
- 581  
582 Elvis Dohmatob, Yunzhen Feng, Pu Yang, Francois Charton, and Julia Kempe. A tale of tails:  
583 Model collapse as a change of scaling laws. In *Forty-first International Conference on Machine*  
584 *Learning*, 2024. URL <https://openreview.net/forum?id=KVvku47shW>.
- 585 Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam  
586 Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for  
587 high-resolution image synthesis. In *Forty-first International Conference on Machine Learning*,  
588 2024.
- 589 Zach Evans, Julian D Parker, CJ Carr, Zack Zukowski, Josiah Taylor, and Jordi Pons. Long-form  
590 music generation with latent diffusion. *arXiv preprint arXiv:2404.10301*, 2024.
- 591  
592 Pierre Fernandez, Guillaume Couairon, Hervé Jégou, Matthijs Douze, and Teddy Furon. The sta-  
593 ble signature: Rooting watermarks in latent diffusion models. In *Proceedings of the IEEE/CVF*  
*International Conference on Computer Vision*, pp. 22466–22477, 2023.

- 594 Shi Fu, Sen Zhang, Yingjie Wang, Xinmei Tian, and Dacheng Tao. Towards theoretical understand-  
595 ings of self-consuming generative models. In *Forty-first International Conference on Machine*  
596 *Learning*, 2024. URL <https://openreview.net/forum?id=aw6L8sB2Ts>.  
597
- 598 Elizabeth Gibney. Ai models fed ai-generated data quickly spew nonsense. *Nature*, 632(8023):  
599 18–19, 2024.
- 600 Josh A Goldstein, Girish Sastry, Micah Musser, Renee DiResta, Matthew Gentzel, and Katerina  
601 Sedova. Generative language models and automated influence operations: Emerging threats and  
602 potential mitigations. *arXiv preprint arXiv:2301.04246*, 2023.  
603
- 604 Xiangming Gu, Chao Du, Tianyu Pang, Chongxuan Li, Min Lin, and Ye Wang. On memorization  
605 in diffusion models. *arXiv preprint arXiv:2310.02664*, 2023.  
606
- 607 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter.  
608 Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in*  
609 *neural information processing systems*, 30, 2017.
- 610 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint*  
611 *arXiv:2207.12598*, 2022.  
612
- 613 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*  
614 *neural information processing systems*, 33:6840–6851, 2020.
- 615 Bernd Jähne. *Digital image processing*. Springer Science & Business Media, 2005.  
616
- 617 Hamidreza Kamkari, Brendan Leigh Ross, Rasa Hosseinzadeh, Jesse C Cresswell, and Gabriel  
618 Loaiza-Ganem. A geometric view of data complexity: Efficient local intrinsic dimension estimat-  
619 ion with diffusion models. In *Thirty-eighth Annual Conference on Neural Information Processing*  
620 *Systems (NeurIPS2024)*, 2024.  
621
- 622 Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-  
623 based generative models. *Advances in Neural Information Processing Systems*, 35:26565–26577,  
624 2022.
- 625 Jae-Eun Lee, Young-Ho Seo, and Dong-Wook Kim. Convolutional neural network-based digital  
626 image watermarking adaptive to the resolution of image and watermark. *Applied Sciences*, 10  
627 (19):6854, 2020.  
628
- 629 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr  
630 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer*  
631 *Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014,*  
632 *Proceedings, Part V 13*, pp. 740–755. Springer, 2014.
- 633 Junxiu Liu, Jiadong Huang, Yuling Luo, Lvchen Cao, Su Yang, Duqu Wei, and Ronglong Zhou.  
634 An optimized image watermarking method based on hd and svd in dwt domain. *IEEE Access*, 7:  
635 80849–80860, 2019.  
636
- 637 Luping Liu, Yi Ren, Zhijie Lin, and Zhou Zhao. Pseudo numerical methods for diffusion models  
638 on manifolds. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=PlKWVd2yBkY>.  
639
- 640 Gabriel Loaiza-Ganem, Brendan Leigh Ross, Rasa Hosseinzadeh, Anthony L. Caterini, and Jesse C.  
641 Cresswell. Deep generative models through the lens of the manifold hypothesis: A survey and  
642 new connections. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL  
643 <https://openreview.net/forum?id=a90WpmSi0I>. Survey Certification, Expert Cer-  
644 tification.  
645
- 646 Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast  
647 ode solver for diffusion probabilistic model sampling in around 10 steps. *Advances in Neural*  
*Information Processing Systems*, 35:5775–5787, 2022a.

- 648 Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. DPM-solver: A  
649 fast ODE solver for diffusion probabilistic model sampling in around 10 steps. In Alice H. Oh,  
650 Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information  
651 Processing Systems*, 2022b. URL [https://openreview.net/forum?id=2uAaGw1P\\_  
652 V](https://openreview.net/forum?id=2uAaGw1P_V).
- 653 Calvin Luo. Understanding diffusion models: A unified perspective. *arXiv preprint  
654 arXiv:2208.11970*, 2022.
- 655 Hila Manor and Tomer Michaeli. On the posterior distribution in denoising: Application to un-  
656 certainty quantification. In *The Twelfth International Conference on Learning Representations*,  
657 2024a. URL [https://openreview.net/forum?id=adSGeugiu\\_j](https://openreview.net/forum?id=adSGeugiu_j).
- 658 Hila Manor and Tomer Michaeli. Zero-shot unsupervised and text-based audio editing using  
659 DDPM inversion. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria  
660 Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st International  
661 Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*,  
662 pp. 34603–34629. PMLR, 21–27 Jul 2024b. URL [https://proceedings.mlr.press/  
663 v235/manor24a.html](https://proceedings.mlr.press/v235/manor24a.html).
- 664 KA Navas, Mathews Cheriyan Ajay, M Lekshmi, Tampy S Archana, and M Sasikumar. Dwt-dct-svd  
665 based watermarking. In *2008 3rd international conference on communication systems software  
666 and middleware and workshops (COMSWARE'08)*, pp. 271–274. IEEE, 2008.
- 667 Weili Nie, Brandon Guo, Yujia Huang, Chaowei Xiao, Arash Vahdat, and Anima Anandkumar.  
668 Diffusion models for adversarial purification. In *International Conference on Machine Learning  
669 (ICML)*, 2022.
- 670 Sandu Popescu, Anthony J Short, and Andreas Winter. Entanglement and the foundations of statis-  
671 tical mechanics. *Nature Physics*, 2(11):754–758, 2006.
- 672 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
673 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
674 models from natural language supervision. In *International conference on machine learning*, pp.  
675 8748–8763. PMLR, 2021.
- 676 Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-  
677 conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022.
- 678 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
679 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-  
680 ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- 681 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar  
682 Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic  
683 text-to-image diffusion models with deep language understanding. *Advances in neural informa-  
684 tion processing systems*, 35:36479–36494, 2022.
- 685 Iliia Shumailov, Zakhar Shumaylov, Yiren Zhao, Nicolas Papernot, Ross Anderson, and Yarin Gal.  
686 Ai models collapse when trained on recursively generated data. *Nature*, 631(8022):755–759,  
687 2024.
- 688 Vassilios Solachidis and Ioannis Pitas. Circularly symmetric watermark embedding in 2-d dft do-  
689 main. *IEEE transactions on image processing*, 10(11):1741–1753, 2001.
- 690 Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Diffusion  
691 art or digital forgery? investigating data replication in diffusion models. In *Proceedings of the  
692 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6048–6058, 2023a.
- 693 Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Under-  
694 standing and mitigating copying in diffusion models. In *Thirty-seventh Conference on Neural  
695 Information Processing Systems*, 2023b. URL [https://openreview.net/forum?id=  
696 HtMXRGbUMt](https://openreview.net/forum?id=HtMXRGbUMt).

- 702 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2021a. URL <https://openreview.net/forum?id=St1giarCHLP>.
- 703  
704  
705
- 706 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben  
707 Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021b. URL <https://openreview.net/forum?id=PXTIG12RRHS>.
- 708  
709
- 710 Jan Pawel Stanczuk, Georgios Batzolis, Teo Deveney, and Carola-Bibiane Schönlieb. Diffu-  
711 sion models encode the intrinsic dimension of data manifolds. In *Forty-first International*  
712 *Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=a0XiA6v256>.
- 713  
714
- 715 Wei Ren Tan, Chee Seng Chan, Hernan Aguirre, and Kiyoshi Tanaka. Improved artgan for condi-  
716 tional synthesis of natural image and artwork. *IEEE Transactions on Image Processing*, 28(1):  
717 394–409, 2019. doi: 10.1109/TIP.2018.2866698. URL <https://doi.org/10.1109/TIP.2018.2866698>.
- 718
- 719 Matthew Tancik, Ben Mildenhall, and Ren Ng. Stegastamp: Invisible hyperlinks in physical pho-  
720 tographs. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recogni-  
721 tion*, pp. 2117–2126, 2020.
- 722
- 723 Peng Wang, Huijie Zhang, Zekai Zhang, Siyi Chen, Yi Ma, and Qing Qu. Diffusion models learn  
724 low-dimensional distributions via subspace clustering. *arXiv preprint arXiv:2409.02426*, 2024.
- 725
- 726 Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment:  
727 from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–  
728 612, 2004.
- 729
- 730 Zijie J. Wang, Evan Montoya, David Munechika, Haoyang Yang, Benjamin Hoover, and  
731 Duen Horng Chau. Large-scale prompt gallery dataset for text-to-image generative models.  
732 *arXiv:2210.14896 [cs]*, 2022. URL <https://arxiv.org/abs/2210.14896>.
- 733
- 734 Yuxin Wen, John Kirchenbauer, Jonas Geiping, and Tom Goldstein. Tree-rings watermarks: Invis-  
735 ible fingerprints for diffusion images. In *Thirty-seventh Conference on Neural Information Pro-  
736 cessing Systems*, 2023a. URL <https://openreview.net/forum?id=Z57JrmubNl>.
- 737
- 738 Yuxin Wen, Yuchen Liu, Chen Chen, and Lingjuan Lyu. Detecting, explaining, and mitigating  
739 memorization in diffusion models. In *The Twelfth International Conference on Learning Repre-  
740 sentations*, 2023b.
- 741
- 742 Zijin Yang, Kai Zeng, Kejiang Chen, Han Fang, Weiming Zhang, and Nenghai Yu. Gaussian shad-  
743 ing: Provable performance-lossless image watermarking for diffusion models. In *Proceedings of  
744 the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12162–12171, 2024.
- 745
- 746 Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and  
747 Yejin Choi. Defending against neural fake news. *Advances in neural information processing  
748 systems*, 32, 2019.
- 749
- 750 Benjamin J Zhang, Siting Liu, Wuchen Li, Markos A Katsoulakis, and Stanley J Osher. Wasserstein  
751 proximal operators describe score-based generative models and resolve memorization. *arXiv  
752 preprint arXiv:2402.06162*, 2024a.
- 753
- 754 Huijie Zhang, Jinfan Zhou, Yifu Lu, Minzhe Guo, Peng Wang, Liyue Shen, and Qing Qu. The  
755 emergence of reproducibility and consistency in diffusion models. In *Forty-first International  
756 Conference on Machine Learning*, 2024b. URL <https://openreview.net/forum?id=HsliOqZkc0>.
- 757
- 758 Kevin Alex Zhang, Lei Xu, Alfredo Cuesta-Infante, and Kalyan Veeramachaneni. Robust invisible  
759 video watermarking with attention. *arXiv preprint arXiv:1909.01285*, 2019.

Lijun Zhang, Xiao Liu, Antoni Viros Martin, Cindy Xiong Bearfield, Yuriy Brun, and Hui Guan. Robust image watermarking using stable diffusion, 2024c. URL <https://arxiv.org/abs/2401.04247>.

Qinsheng Zhang and Yongxin Chen. Fast sampling of diffusion models with exponential integrator. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=Loek7hfb46P>.

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.

Wenliang Zhao, Lujia Bai, Yongming Rao, Jie Zhou, and Jiwen Lu. Unipc: A unified predictor-corrector framework for fast sampling of diffusion models. In *NeurIPS*, 2023a. URL [http://papers.nips.cc/paper\\_files/paper/2023/hash/9c2aa1e456ea543997f6927295196381-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/9c2aa1e456ea543997f6927295196381-Abstract-Conference.html).

Xuandong Zhao, Kexun Zhang, Yu-Xiang Wang, and Lei Li. Generative autoencoders as watermark attackers: Analyses of vulnerabilities and threats. 2023b.

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 as:

$$[\mathbf{x}_0^{*\mathcal{W}}]_i = \text{ChannelAverage}(\mathbf{x}_0^{\mathcal{W}}, \mathbf{x}_0^*), \tag{11}$$

$$= \begin{cases} [\mathbf{x}_0^{\mathcal{W}}]_i, i = c \\ (1 - \gamma)[\mathbf{x}_0^{\mathcal{W}}]_i + \gamma[\mathbf{x}_0^*]_i, i \neq c \end{cases} \tag{12}$$

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.4 for a more detailed analysis.

## C ADDITIONAL EXPERIMENTS

### C.1 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 (JPEG) with a compression rate of 25%.
- Gaussian blurring (G.Blur) with an  $8 \times 8$  filter size.
- Gaussian noise (G.Noise) with  $\sigma = 0.1$ .
- Color jitter (CJ) with brightness factor uniformly ranges between 0 and 6.

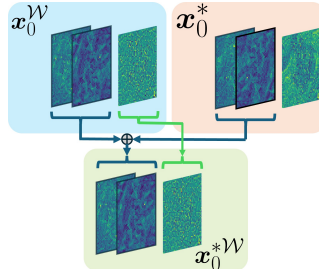


Figure 6: Illustration of channel average



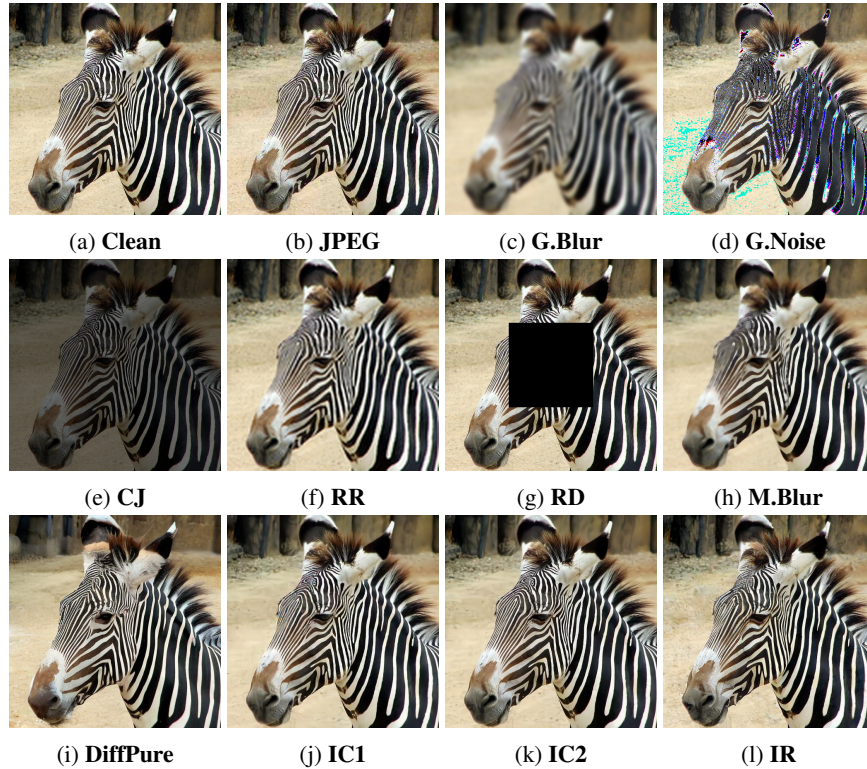


Figure 7: Visualization of different attacks.

- Resize and restore (RR). Resize to 50% of pixels and restore to original size.
- Random drop (RD). Random drop a square with 40% of pixels.
- Median blurring (M.Blur) with a  $7 \times 7$  median filter.
- Diffusion purification Nie et al. (2022) (DiffPure) with the purified step at 0.3T.
- VAE-based image compression Cheng et al. (2020) (IC1) and Ballé et al. (2018) (IC2), with a quality level of 3.
- Diffusion-based image regeneration Zhao et al. (2023b) with 60 denoising steps.

Visualizations of these attacks are in Figure 7. Detailed experiments for table 1 (Table 4) on above attacks are in Table 3 (Table 4).

Table 3: Watermarking Robustness for different attacks under the server scenario.

Method	Watermarking Robustness (AUC ↑/TPR@1%FPR↓)												
	Clean	JPEG	G.Blur	G.Noise	CJ	RR	RD	M.Blur	DiffPure	IC1	IC2	IR	Average
DwtDet	0.97/0.85	0.47/0.00	0.51/0.02	0.96/0.78	0.53/0.15	0.66/0.14	0.99/0.88	0.58/0.01	0.50/0.00	0.52/0.01	0.49/0.00	0.50/0.00	0.61/0.18
DwtDetSvd	<b>1.00/1.00</b>	0.64/0.10	0.96/0.70	0.99/0.99	0.53/0.12	0.99/0.99	1.00/1.00	<b>1.00/1.00</b>	0.51/0.02	0.73/0.03	0.68/0.04	0.70/0.07	0.79/0.46
RivaGAN	1.00/0.99	0.94/0.69	0.96/0.76	0.97/0.88	0.95/0.79	0.99/0.98	0.99/0.98	0.99/0.97	0.73/0.16	0.65/0.03	0.63/0.04	0.56/0.00	0.85/0.57
Stegastamp	<b>1.00/1.00</b>	<b>1.00/1.00</b>	1.00/0.95	0.98/0.97	1.00/0.97	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	0.81/0.29	1.00/0.97	1.00/0.99	0.90/0.43	0.97/0.87
Stable Signature	<b>1.00/1.00</b>	0.99/0.76	0.57/0.00	0.71/0.14	0.96/0.87	0.90/0.34	<b>1.00/1.00</b>	0.95/0.62	0.54/0.01	0.93/0.58	0.91/0.50	0.67/0.02	0.83/0.44
Tree-Ring Watermarks	<b>1.00/1.00</b>	0.99/0.97	0.98/0.98	0.94/0.50	0.96/0.67	<b>1.00/1.00</b>	0.99/0.97	0.99/0.94	0.98/0.73	0.99/0.97	0.99/0.98	0.99/0.92	0.98/0.87
RingID	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	1.00/0.99	0.99/0.98	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>
Gaussian Shading	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>
Shallow Diffuse (ours)	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>

## C.2 MULTI-KEY WATERMARKING

In this section, we examine the capability of Shallow Diffuse to support multi-key watermarking. We evaluate two important tasks associated with multi-key watermarking: Multi-key identification and Multi-key re-watermarking.

Table 4: Watermarking Robustness for different attacks under the user scenario.

Method	Watermarking Robustness (AUC ?/TPR@1%FPR)												Average
	Clean	JPEG	G.Blur	G.Noise	CJ	RR	RD	M.Blur	DiffPure	IC1	IC2	IR	
<b>COCO Dataset</b>													
DwtDct	0.98/0.83	0.50/0.01	0.50/0.00	0.97/0.81	0.54/0.14	0.67/0.17	0.99/0.93	0.59/0.05	0.46/0.00	0.49/0.00	0.49/0.01	0.46/0.00	0.61/0.19
DwtDctSvd	<b>1.00/1.00</b>	0.64/0.13	0.98/0.83	0.99/0.99	0.54/0.13	1.00/1.00	1.00/1.00	<b>1.00/1.00</b>	0.50/0.01	0.70/0.05	0.64/0.04	0.68/0.07	0.79/0.48
RivaGAN	<b>1.00/1.00</b>	0.97/0.86	0.98/0.86	0.99/0.94	0.96/0.82	1.00/1.00	1.00/1.00	1.00/1.00	0.63/0.02	0.68/0.05	0.66/0.04	0.75/0.15	0.87/0.61
Stegastamp	<b>1.00/1.00</b>	<b>1.00/1.00</b>	0.99/0.90	0.90/0.87	1.00/0.98	1.00/0.99	1.00/0.99	<b>1.00/1.00</b>	0.81/0.27	1.00/0.95	1.00/0.95	0.85/0.28	0.96/0.83
Tree-Ring Watermarks	<b>1.00/1.00</b>	0.99/0.87	0.99/0.86	<b>1.00/1.00</b>	0.88/0.49	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	0.99/0.93
RingID	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	0.98/0.86	<b>1.00/0.99</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/0.99</b>
Shallow Diffuse (ours)	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/0.99</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	0.99/0.86	1.00/0.99	0.99/0.97	<b>1.00/1.00</b>	1.00/0.98
<b>DiffusionDB Dataset</b>													
DwtDct	0.96/0.76	0.47/0.002	0.51/0.018	0.96/0.78	0.53/0.15	0.66/0.14	0.99/0.88	0.58/0.01	0.50/0.004	0.52/0.008	0.49/0.004	0.50/0.002	0.61/0.18
DwtDctSvd	<b>1.00/1.00</b>	0.64/0.10	0.96/0.70	0.99/0.99	0.53/0.12	1.00/1.00	1.00/1.00	<b>1.00/1.00</b>	0.51/0.022	0.73/0.03	0.68/0.04	0.70/0.07	0.79/0.46
RivaGAN	1.00/0.98	0.94/0.69	0.96/0.76	0.97/0.88	0.95/0.79	1.00/0.98	0.99/0.98	<b>1.00/1.00</b>	0.56/0.004	0.65/0.03	0.63/0.04	0.73/0.16	0.85/0.57
Stegastamp	<b>1.00/1.00</b>	<b>1.00/1.00</b>	0.99/0.88	0.91/0.89	1.00/0.99	1.00/0.97	1.00/1.00	1.00/0.96	0.83/0.28	1.00/0.91	1.00/0.93	0.85/0.40	0.96/0.84
Tree-Ring Watermarks	<b>1.00/1.00</b>	0.99/0.68	0.94/0.62	<b>1.00/1.00</b>	0.84/0.15	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	0.99/0.99	0.99/0.99	0.99/0.98	0.96/0.92	0.97/0.85
RingID	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	0.98/0.86	1.00/0.98	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/0.99</b>
Shallow Diffuse (ours)	<b>1.00/1.00</b>	1.00/0.99	1.00/0.99	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	0.96/0.90	0.96/0.92	0.97/0.93	0.98/0.96	0.99/0.97
<b>WikiArt Dataset</b>													
DwtDct	0.96/0.75	0.46/0.004	0.51/0.008	0.95/0.75	0.50/0.13	0.68/0.13	0.98/0.87	0.61/0.08	0.48/0.006	0.47/0.006	0.49/0.002	0.48/0.006	0.60/0.18
DwtDctSvd	<b>1.00/1.00</b>	0.65/0.22	0.97/0.76	0.99/0.99	0.50/0.10	1.00/1.00	1.00/1.00	<b>1.00/1.00</b>	0.47/0.03	0.72/0.04	0.66/0.07	0.67/0.08	0.78/0.48
RivaGAN	<b>1.00/1.00</b>	0.96/0.80	0.99/0.95	0.98/0.93	0.89/0.66	1.00/1.00	1.00/1.00	<b>1.00/1.00</b>	0.63/0.02	0.66/0.04	0.67/0.04	0.80/0.11	0.87/0.60
Stegastamp	<b>1.00/1.00</b>	1.00/0.96	0.97/0.77	0.92/0.88	0.98/0.84	0.99/0.89	<b>1.00/1.00</b>	0.99/0.91	0.77/0.20	0.99/0.95	0.99/0.90	0.80/0.30	0.95/0.75
Tree-Ring Watermarks	<b>1.00/1.00</b>	1.00/0.97	1.00/0.88	<b>1.00/1.00</b>	0.71/0.26	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	0.97/0.92
RingID	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	0.95/0.82	0.99/0.98	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	0.99/0.98
Shallow Diffuse (ours)	<b>1.00/1.00</b>	1.00/0.99	1.00/0.99	<b>1.00/1.00</b>	<b>1.00/0.99</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	0.97/0.94	0.98/0.95	<b>1.00/1.00</b>	<b>1.00/0.99</b>

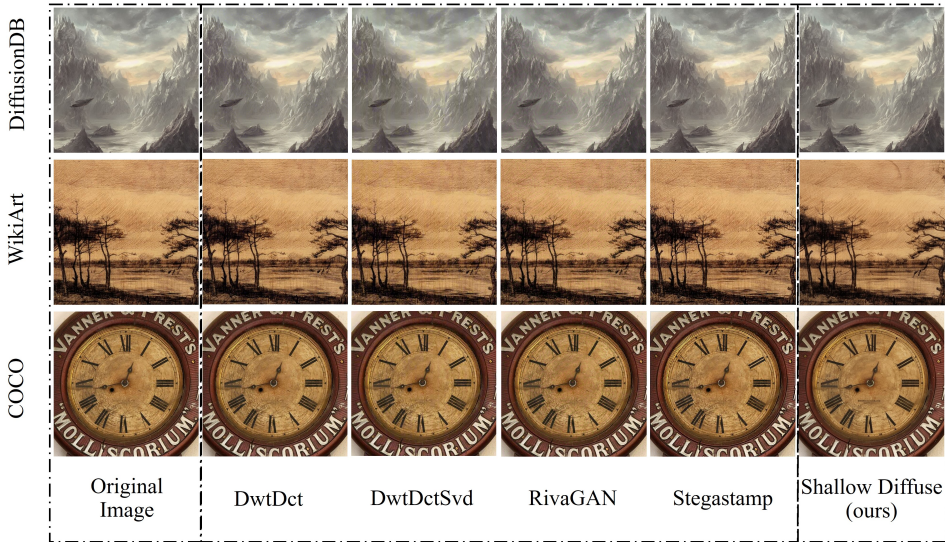


Figure 8: **Generation Consistency in User Scenarios.** We compare the visualization quality of our method against DwtDct, DwtDctSvd, RivaGAN, and Stegastamp across the DiffusionDB, WikiArt, and COCO datasets.

**Multi-key identification** This is a classification task designed to test the ability to accurately identify individual watermarks. We generate a set of  $N = 2048$  watermarks, all using the same circular mask  $M$  but with distinct ring-shaped keys  $\{W_i\}_{i=1}^N$ . During watermarking, a random key  $W_i$  is selected and injected into images. After an attack is applied, we attempt to detect the watermark key  $W_j$  and determine if  $i = j$ . The success rate of identification serves as the evaluation metric. This setup is inspired by the work in Ci et al. (2024). We compare Tree-Ring, RingID, and Shallow Diffuse in the server scenario. The results of this experiment are shown in Table 5. Despite lacking a dedicated design for multi-key scenarios, Shallow Diffuse outperforms Tree-Ring, RingID, specifically designed for multi-key identification, achieves the highest success rate. Exploring multi-key identification strategies could be an important direction for future research.

**Multi-key re-watermarking** : This task evaluates the ability to embed multiple watermarks into the same image and detect each one independently. For this experiment, we test cases with 2, 4, 8,

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025



Figure 9: **Generation Consistency in server scenarios.** We compare the visualization quality of our method against the original image and StageStamp.

16, 32 watermarks. Each watermark uses a unique ring-shaped key  $W_i$  and a non-overlapped mask  $M$  (part of a circle). This is a non-trivial setting as we could pre-defined the key number and non-overlapped mask  $M$  for application. The metric for this task is the average robustness across all keys, measured in terms of AUC and TPR@1%FPR. For this study, we test the Tree-Ring and Shallow Diffuse in the server scenario. The results of this experiment are presented in Table 6. Shallow Diffuse consistently outperformed Tree-Ring in robustness across different numbers of users. Even as the number of users increased to 32, Shallow Diffuse maintained strong robustness under clean conditions. However, in adversarial settings, its robustness began to decline when the number of users exceeded 16. Under the current setup, when the number of users surpasses the predefined limit, our method becomes less robust and accurate. We believe that enabling watermarking for hundreds or even thousands of users simultaneously is a challenging yet promising future direction for Shallow Diffuse.

Table 5: **Multi-key identification for different attacks under the server scenario.**

Method	Successfull Rate $\uparrow$									
	Clean	JPEG	G.Blur	G.Noise	CJ	DiffPure	IC1	IC2	IR	Average
Tree-Ring	0.20	0.04	0.09	0.07	0.06	0.06	0.28	0.29	0.23	0.15
RingID	<b>1.00</b>	<b>0.97</b>	<b>0.97</b>	<b>0.95</b>	<b>0.87</b>	<b>0.88</b>	<b>0.98</b>	<b>0.98</b>	<b>0.99</b>	<b>0.95</b>
Shallow Diffuse	0.88	0.77	0.57	0.88	0.40	0.48	0.41	0.64	0.80	0.65

### C.3 ABLATION STUDY OF DIFFERENT WATERMARK PATTERNS

In Table 7, we examine various combinations of watermark patterns  $M \odot W$ . For the shape of the mask  $M$ , "Circle" refers to a circular mask  $M$  (see Figure 2 top left), while "Ring" represents a ring-shaped  $M$ . Since the mask is centered in the middle of the figure, "Low" and "High" denote frequency regions: "Low" represents a DFT with zero-frequency centering, whereas "High"

Table 6: Multi-key re-watermark for different attacks under the server scenario.

Watermark number	Method	Watermarking Robustness (AUC $\uparrow$ /TPR@1%FPR $\uparrow$ )												
		Clean	JPEG	G.Blur	G.Noise	CJ	RR	RD	M.Blur	DiffPure	IC1	IC2	IR	Average
2	Tree-Ring	<b>1.00/1.00</b>	0.99/0.84	1.00/0.97	0.95/0.83	0.98/0.75	1.00/1.00	1.00/1.00	1.00/1.00	0.91/0.23	1.00/0.91	0.98/0.82	0.94/0.49	0.98/0.80
	Shallow Diffuse	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>0.98/0.95</b>	<b>1.00/0.90</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>0.98/0.65</b>	<b>1.00/0.91</b>	<b>1.00/0.97</b>	<b>1.00/0.99</b>	<b>0.99/0.95</b>
4	Tree-Ring	<b>1.00/1.00</b>	0.98/0.63	1.00/0.89	0.96/0.86	0.90/0.54	1.00/0.92	1.00/0.99	1.00/0.95	0.88/0.11	0.99/0.72	0.97/0.67	0.92/0.37	0.96/0.70
	Shallow Diffuse	<b>1.00/1.00</b>	1.00/0.96	0.99/0.88	0.97/0.91	0.99/0.82	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	0.94/0.37	0.99/0.80	0.99/0.83	0.99/0.89	0.99/0.86
8	Tree-Ring	1.00/0.95	0.90/0.32	0.97/0.56	0.92/0.64	0.90/0.45	0.98/0.71	1.00/0.89	0.98/0.68	0.77/0.08	0.91/0.38	0.89/0.25	0.83/0.16	0.91/0.47
	Shallow Diffuse	<b>1.00/1.00</b>	0.99/0.85	0.97/0.73	0.97/0.90	0.98/0.80	1.00/0.98	<b>1.00/1.00</b>	1.00/0.96	0.91/0.36	0.98/0.71	0.97/0.70	0.99/0.80	0.98/0.80
16	Tree-Ring	0.96/0.57	0.78/0.18	0.87/0.32	0.87/0.38	0.84/0.24	0.90/0.42	0.95/0.53	0.90/0.36	0.68/0.05	0.80/0.18	0.77/0.14	0.72/0.05	0.83/0.26
	Shallow Diffuse	1.00/0.89	0.94/0.59	0.89/0.39	0.94/0.73	0.92/0.53	0.97/0.73	0.99/0.84	0.96/0.73	0.78/0.11	0.90/0.46	0.91/0.46	0.92/0.55	0.92/0.56
32	Tree-Ring	0.95/0.44	0.77/0.11	0.85/0.15	0.86/0.31	0.80/0.15	0.88/0.22	0.94/0.34	0.89/0.26	0.63/0.03	0.78/0.11	0.75/0.08	0.70/0.05	0.80/0.16
	Shallow Diffuse	0.99/0.89	0.91/0.46	0.86/0.26	0.93/0.63	0.91/0.47	0.96/0.65	0.99/0.84	0.95/0.59	0.74/0.07	0.87/0.31	0.87/0.30	0.89/0.28	0.90/0.44

indicates a DFT without zero-frequency centering, as illustrated in Figure 2 bottom. For the distribution of  $\mathbf{W}$ , "Zero" implies all values are zero, "Rand" denotes values sampled from  $\mathcal{N}(0, 1)$ , and "Rotational Rand" represents multiple concentric rings in  $\mathbf{W}$ , with each ring's values sampled from  $\mathcal{N}(0, 1)$ .

As shown in Table 7, watermarking in high-frequency regions (Rows 7-9) yields improved image consistency compared to low-frequency regions (Rows 1-6). Additionally, the "Circle"  $M$  combined with "Rotational Rand"  $\mathbf{W}$  (Rows 3 and 9) demonstrates greater robustness than other watermark patterns. Consequently, Shallow Diffuse employs the "Circle"  $M$  with "Rotational Rand"  $\mathbf{W}$  in the high-frequency region.

Table 7: Ablation study on different watermark patterns.

Method & Dataset			PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	Average Watermarking Robustness (AUC $\uparrow$ /TPR@1%FPR $\uparrow$ )
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	<b>1.00/1.00</b>
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	<b>42.3</b>	<b>0.98</b>	<b>0.004</b>	0.86/0.35
High	Circle	Rotational Rand	38.0	0.94	0.01	<b>1.00/1.00</b>

#### C.4 ABLATION STUDY OF CHANNEL AVERAGE

We evaluate Shallow Diffuse with channel averaging enabled ( $\gamma = 1.0$ ) and disabled ( $\gamma = 0.0$ ), as shown in Table 8. Unlike the adaptive image enhancement techniques proposed in Zhang et al. (2024c), our approach embeds the watermark in a single channel while averaging the non-watermarked channels. This design takes advantage of the fact that many image processing operations, such as color jittering or Gaussian blurring, tend to affect all channels uniformly. By isolating the watermark in one channel, it will be less vulnerable to those attacks. Thus, applying channel averaging slightly enhances robustness against certain attacks while maintaining comparable consistency. Therefore, we set  $\gamma = 1.0$  for Shallow Diffuse.

Table 8: ablation study on channel average.

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

#### C.5 ABLATION STUDY OF WATERMARKING EMBEDDED CHANNEL.

As shown in Table 9, we evaluate specific embedding channels  $c$  for Shallow Diffuse, where "0," "1," "2," and "3" denote  $c = 0, 1, 2, 3$ , respectively, and "0 + 1 + 2 + 3" indicates watermarking applied across all channels<sup>2</sup>. Since applying watermarking to any single channel yields similar results (Row

<sup>2</sup>Here we apply Shallow Diffuse on the latent space of Stable Diffusion, the channel dimension is 4.

1-4), but applying it to all channels (Row 5) negatively impacts image consistency and robustness, we set  $c = 3$  for Shallow Diffuse. This finding aligns with the observations in the channel average ablation study (appendix C.4). The reason is that many image processing operations tend to affect all channels uniformly, making watermarking across all channels more susceptible to such attacks.)

Table 9: Ablation study on 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
0	36.46	<b>0.93</b>	<b>0.02</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	0.99
1	36.57	<b>0.93</b>	<b>0.02</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	0.99
2	36.13	0.92	<b>0.02</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>
3	<b>36.64</b>	<b>0.93</b>	<b>0.02</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>
0 + 1 + 2 + 3	33.19	0.83	0.05	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	0.95

## C.6 ABLATION STUDY OF DIFFERENT SAMPLING METHODS

We conducted ablation studies on various diffusion model sampling methods, including DDIM, DEIS Zhang & Chen (2023), DPM-Solver Lu et al. (2022b), PNDM Liu et al. (2022), and UniPC Zhao et al. (2023a). All methods were evaluated using 50 sampling steps. The results, presented in Table 10, indicate that Shallow Diffuse is not highly sensitive to the choice of sampling method. Across all methods, the generation quality and watermark robustness remain consistent.

Table 10: Ablation study on sampling methods.

Sampling Method	Generation Quality	Watermark Robustness (AUC $\uparrow$ /TPR@1%FPR $\uparrow$ )					
	CLIP-Score $\uparrow$	Clean	JPEG	G.Blur	G.Noise	CJ	Adversarial Average
DDIM	<b>0.3652</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>
DEIS	0.3651	<b>1.00/1.00</b>	0.99/0.99	<b>1.00/1.00</b>	<b>1.00/1.00</b>	0.99/0.95	1.00/0.99
DPM-Solver	0.3645	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	1.00/0.99	0.99/0.94	1.00/0.98
PNDM	0.3651	<b>1.00/1.00</b>	0.99/0.99	<b>1.00/1.00</b>	<b>1.00/1.00</b>	0.98/0.96	1.00/0.99
UniPC	0.3645	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>	<b>1.00/1.00</b>

## D PROOFS IN SECTION 4

### D.1 PROOFS OF THEOREM 1

*Proof of Theorem 1.* According to Assumption 1, we have  $\|\hat{\mathbf{x}}_{0,t}^{\mathcal{V}} - \hat{\mathbf{x}}_{0,t}\|_2^2 = \lambda \|\mathbf{J}_{\theta,t}(\mathbf{x}_t) \cdot \Delta \mathbf{x}\|_2^2$ . 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:

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

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  $\epsilon$ :

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

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:

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

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  $\|\mathbf{J}_{\boldsymbol{\theta}, t-1}(\mathbf{x}_{t-1}) \mathbf{J}_{\boldsymbol{\theta}, t}(\mathbf{x}_t) \Delta \mathbf{x}\|_2 \leq \|\mathbf{J}_{\boldsymbol{\theta}, t-1}(\mathbf{x}_{t-1})\|_2 \|\mathbf{J}_{\boldsymbol{\theta}, t}(\mathbf{x}_t) \Delta \mathbf{x}\|_2 \leq L \|\mathbf{J}_{\boldsymbol{\theta}, t}(\mathbf{x}_t) \Delta \mathbf{x}\|_2$ . 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  $1 - \frac{1}{r_t} - \frac{1}{r_{t-1}}$ ,

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

□

## E 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 vectors  $\{v_1, v_2, \dots, v_i, \dots, v_d\}$  are orthonormal 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:

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

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

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

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

Additionally,  $\sum_{j=1}^d Y_j = 1$ . So:

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

□

**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:

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

□