

Forging and Removing Latent-Noise Diffusion Watermarks

Using a Single Image

Anubhav Jain¹, Yuya Kobayash², Naoki Murata², Yuhta Takida², Takashi Shibuya²,
Yuki Mitsufuji^{2,3}, Niv Cohen^{1,*}, Nasir Memon^{1,*}, Julian Togelius^{1,*}

¹New York University, ²Sony AI, ³Sony Group Corporation

arXiv preprint
arXiv:2504.20111

Presenter: Yufei Zhang
2026.03.22

Introduction

Visible Watermark



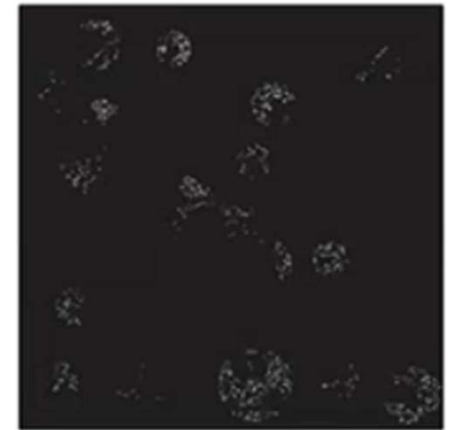
Invisible Watermark



Clean



Watermarked



- Least Significant Bit (Wolfgang et al., 1996), Spatial Domain (Ghazanfari et al., 2011), Frequency Domain (Holub et al., 2012)
- DL-based approaches: SteganoGAN (Zhang et al., 2019a), Self-Supervised Learning (Fernandez et al., 2021)

Introduction



Emergence of invisible watermarking in AI industry...

Google DeepMind ▾ Models Research Science About Build with Gemini

Overview How it works **Hands-on**

Detecting SynthID watermarks in Gemini
Want to check if an image, video or audio clip was generated, or edited, by Google AI? Ask Gemini.

Simply upload the image, video or audio clip ask if it's been created or altered by Google check for a SynthID watermark, and let you know

[Try in Gemini >](#)

[How to verify in Gemini >](#)

Meta AI Research The Latest About Get Llama

Try Meta AI



Research

Stable Signature: A new method for watermarking images created by open source generative AI

October 6, 2023 · 6 minute read

RESEARCH

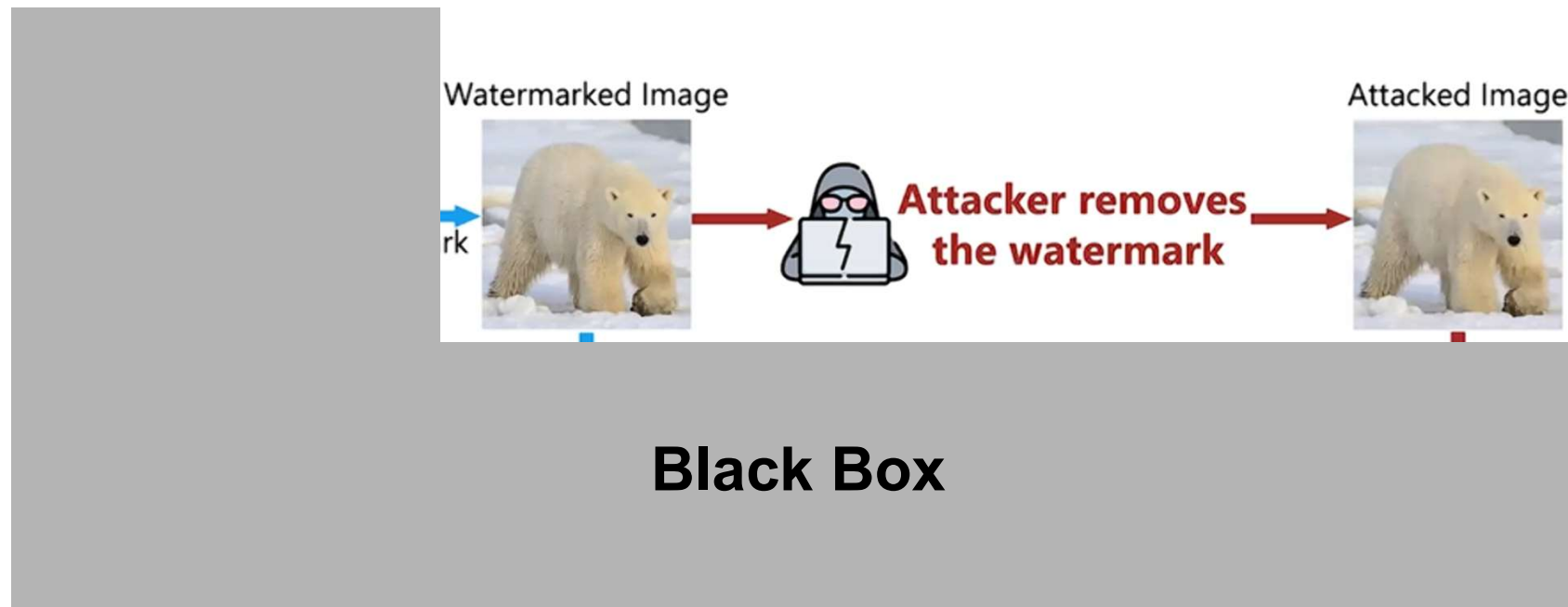
Introducing a watermarking method to distinguish images created by Generative AI

Meta

AI at Meta

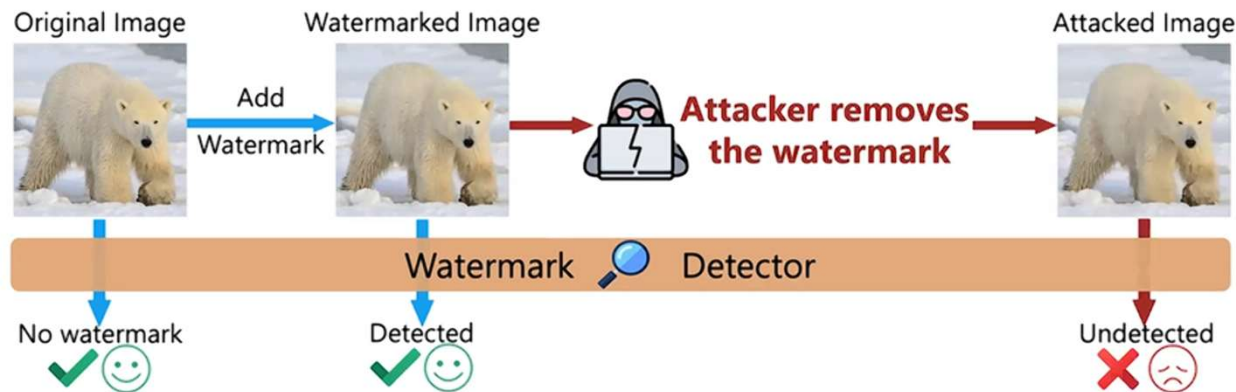
Introduction

Are Invisible Watermarks Good enough?



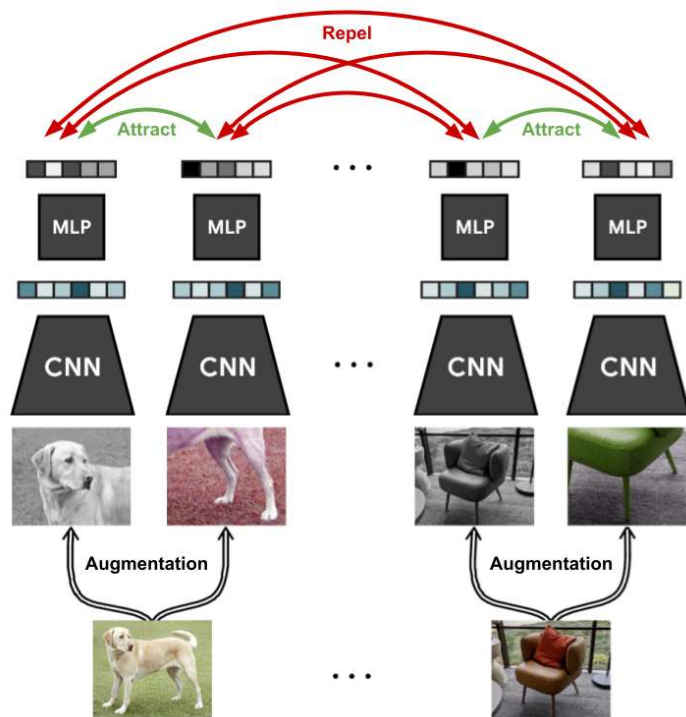
Introduction : Watermarking Schemes

Are Invisible Watermarks Good enough?

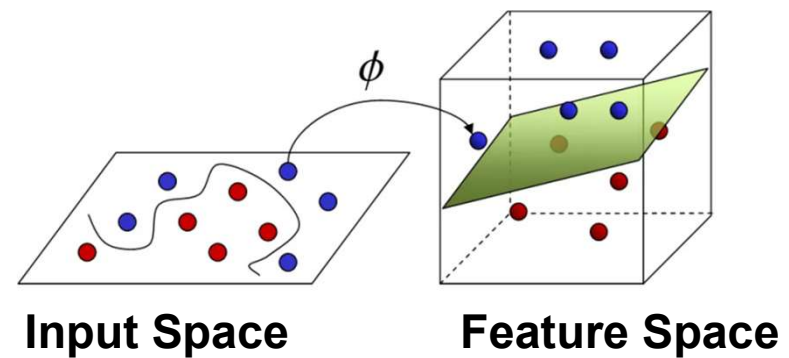


Background : Watermarking Schemes

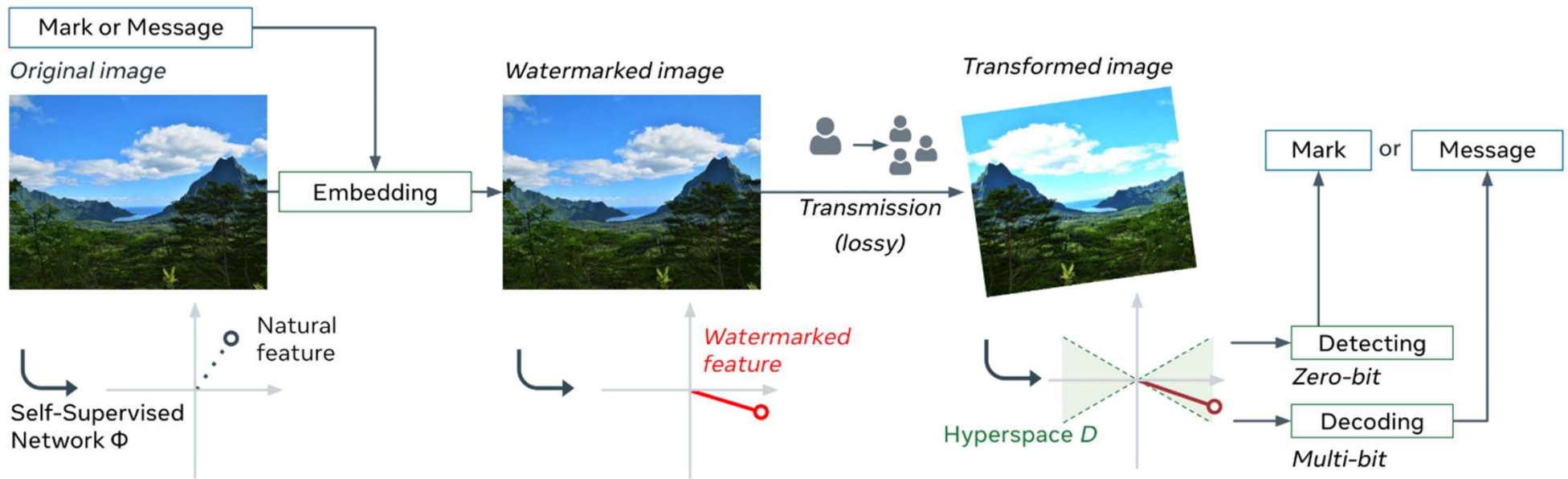
Self-Supervised Learning (SSL)



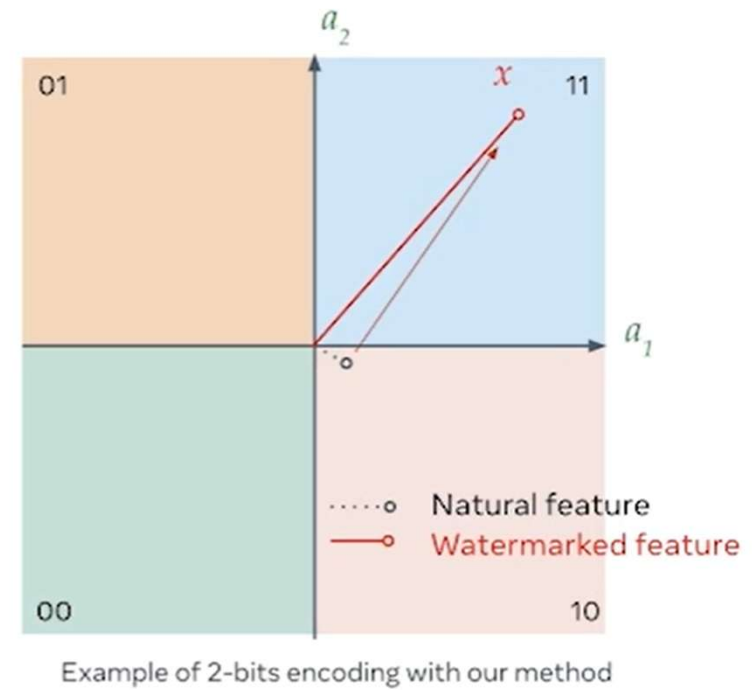
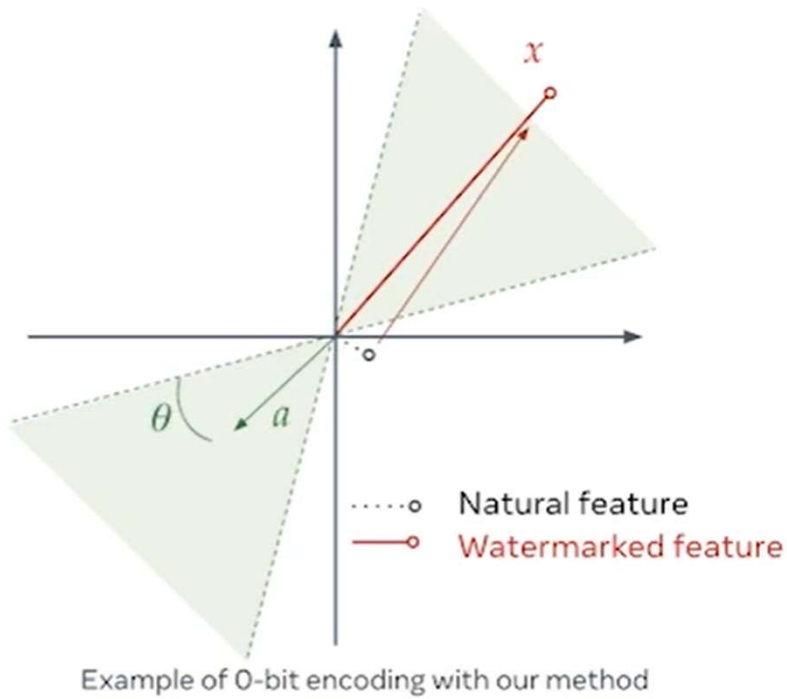
- 1. Intrinsic robustness of SSL to image transformations
- 2. Does not suffer from semantic collapse of supervised learning (More capacity for Watermarking)



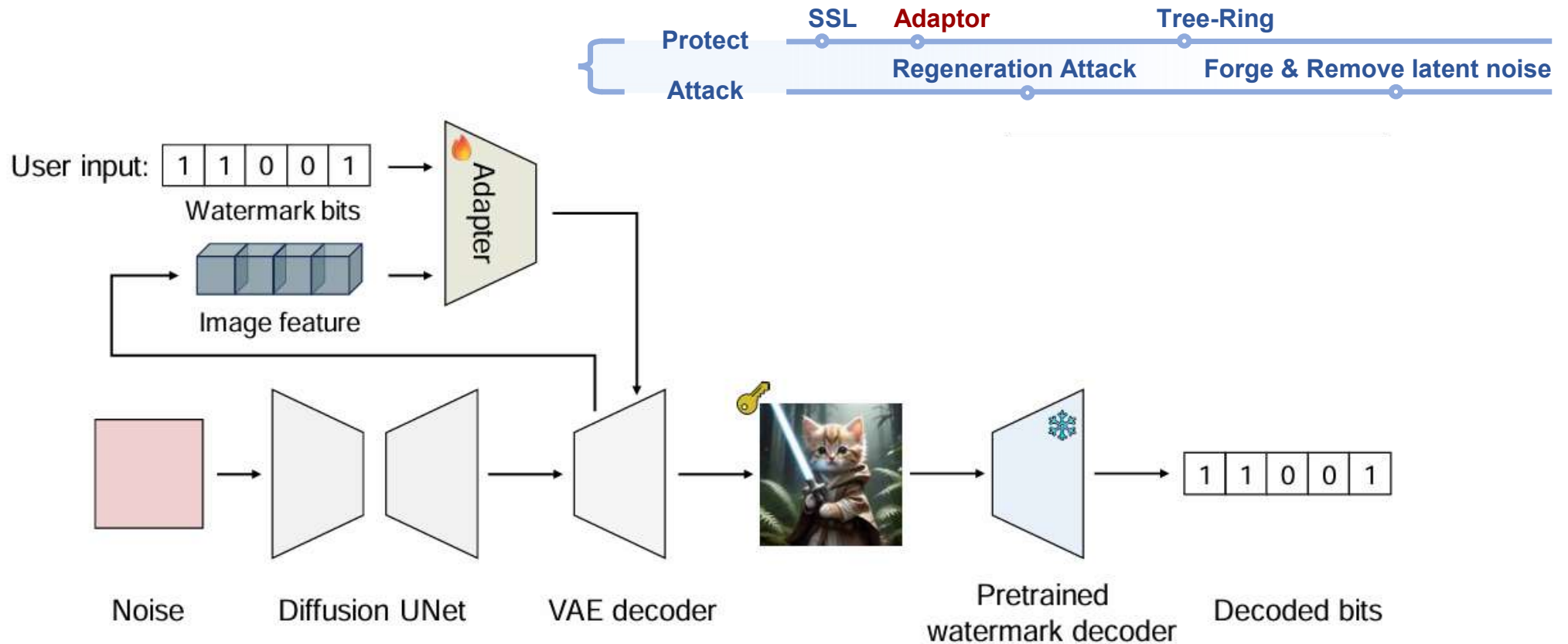
Background : Watermarking Schemes



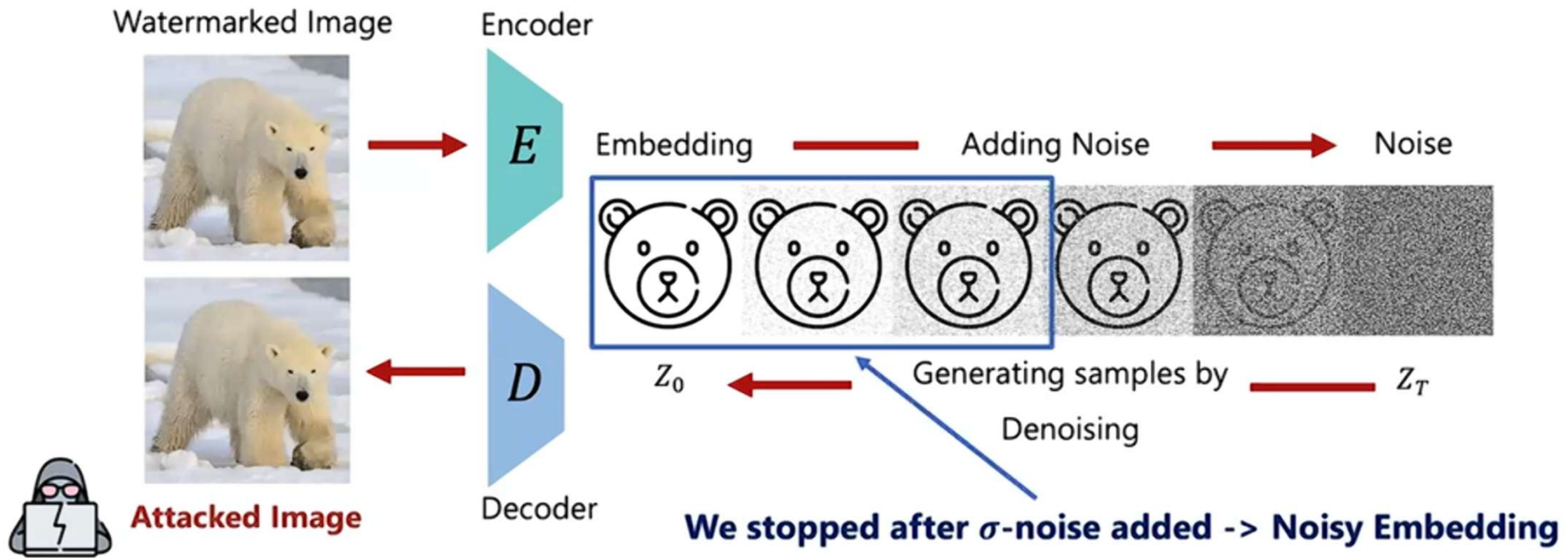
Background : Watermarking Schemes



Background : Watermarking Schemes



Background : Watermarking Schemes



Background : Watermarking Schemes

Guaranteed to remove any Pixelbased Invisible Watermarks.

MS-COCO Dataset					
Attacker	PSNR↑	SSIM↑	FID↓	Bit Acc↓	Detect Acc↓
Stable Signature watermarking:					
Brightness 0.5	28.53	0.864	11.75	0.967	0.990
Contrast 0.5	27.20	0.842	10.73	0.965	0.990
JPEG 50	29.37	0.873	15.01	0.866	0.966
Gaussian noise	25.46	0.788	30.60	0.920	0.972
Gaussian blur	25.48	0.798	17.72	0.896	0.986
BM3D denoise	29.24	0.871	31.65	0.946	0.952
VAE-Bmshj2018	28.95	0.867	31.86	0.636	0.248
VAE-Cheng2020	28.67	0.864	29.43	0.682	0.442
Diffusion model	29.33	0.879	20.64	0.486	0.000

Removing 100% Stable Signature Watermarks from Meta AI via diffusion attacks.

Background : Watermarking Schemes

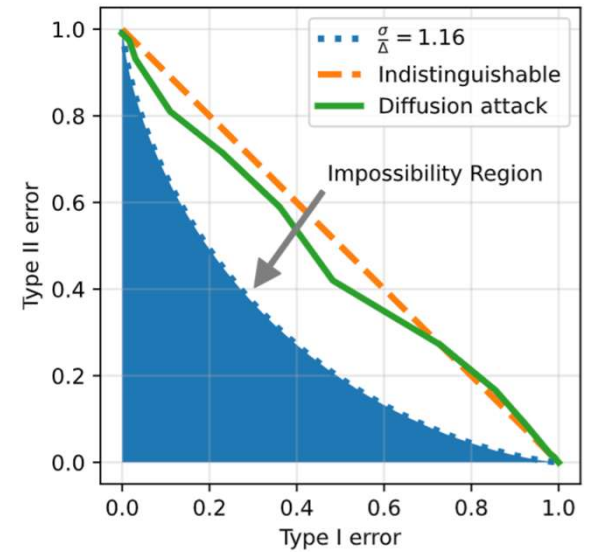
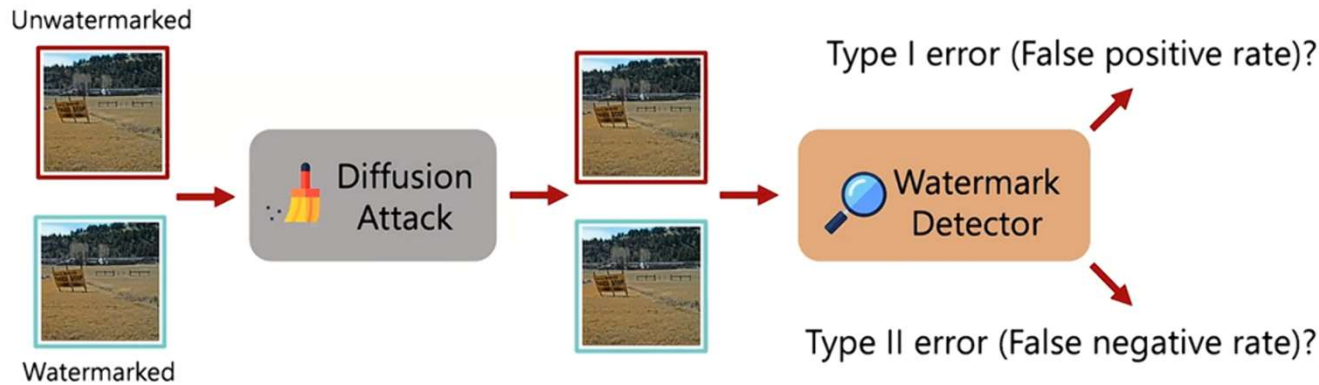
Definition 1 (Δ -Invisible watermark).

$$\text{dist}(x, x_w) \leq \Delta. \quad \|x - x_w\| \leq \Delta.$$

Definition 2 (f-Certified-Watermark-Free).

$$\varepsilon_2 \geq f(\varepsilon_1), \quad \forall \varepsilon_1 \in [0,1].$$

x	clean
x_w	watermarked
ε_1	FP rate
ε_2	FN rate



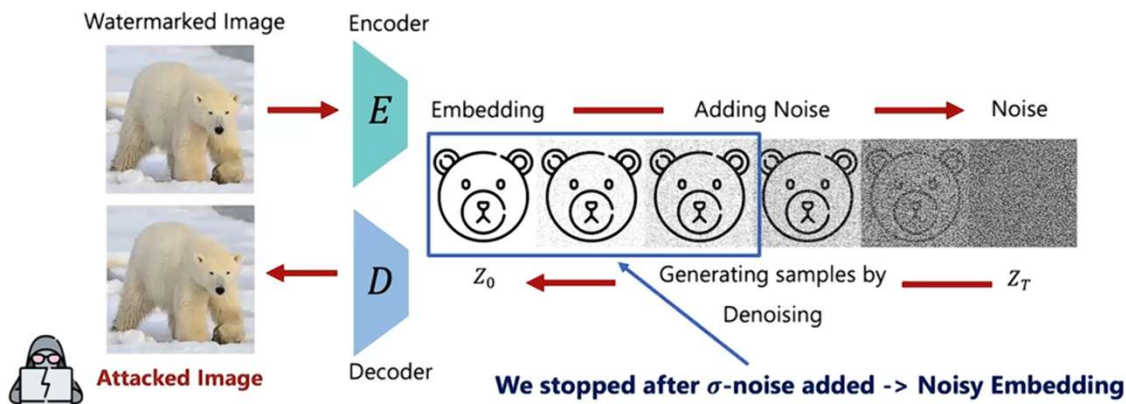
Background : Watermarking Schemes

Definition 1 (Δ -Invisible watermark). $\text{dist}(x, x_w) \leq \Delta. \quad \|x - x_w\| \leq \Delta.$

Definition 2 (f-Certified-Watermark-Free). $\varepsilon_2 \geq f(\varepsilon_1), \quad \forall \varepsilon_1 \in [0,1].$

x	clean
x_w	watermarked
ε_1	FP rate
ε_2	FN rate

Definition 3 (Local Watermark-Specific Lipschitz property). $|\phi(x_w) - \phi(x)| \leq L_{x,w}|x_w - x|$



$$\begin{cases} \tilde{x} \sim A(\phi(x_0) + N(0, \sigma^2 I)) \\ \tilde{x} \sim A(\phi(x_w) + N(0, \sigma^2 I)) \end{cases}$$

Background : Watermarking Schemes



Definition 1 (Δ -Invisible watermark). $\text{dist}(x, x_w) \leq \Delta. \quad \|x - x_w\| \leq \Delta.$

Definition 2 (f-Certified-Watermark-Free). $\varepsilon_2 \geq f(\varepsilon_1), \quad \forall \varepsilon_1 \in [0,1].$

x	clean
x_w	watermarked
ε_1	FP rate
ε_2	FN rate

Definition 3 (Local Watermark-Specific Lipschitz property). $|\phi(x_w) - \phi(x)| \leq L_{x,w}|x_w - x|$

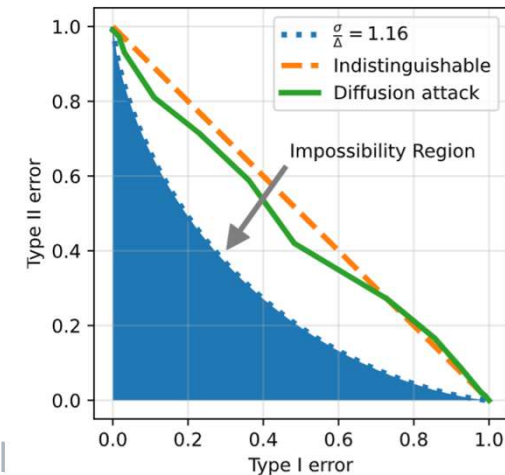
$$\left\{ \begin{array}{l} \tilde{x} \sim A(\phi(x_0) + N(0, \sigma^2 I)) \\ \tilde{x} \sim A(\phi(x_w) + N(0, \sigma^2 I)) \end{array} \right\} \left\{ \begin{array}{l} N(\phi(x_0), \sigma^2 I) \\ N(\phi(x_w), \sigma^2 I) \end{array} \right\} \xrightarrow{\begin{array}{l} \mu = \phi(x_w) - \phi(x_0) \\ \|\mu\| \leq L_{x,w}\Delta \end{array}} \left\{ \begin{array}{l} N(0, 1) \\ N\left(\frac{L_{x,w}\Delta}{\sigma}, 1\right) \end{array} \right\}$$

Theorem1 (Neyman-Pearson Lemma).

The likelihood ratio test is uniform most powerful.

Theorem2 (Gaussian Differential Privacy [Dong et al. 2022]).

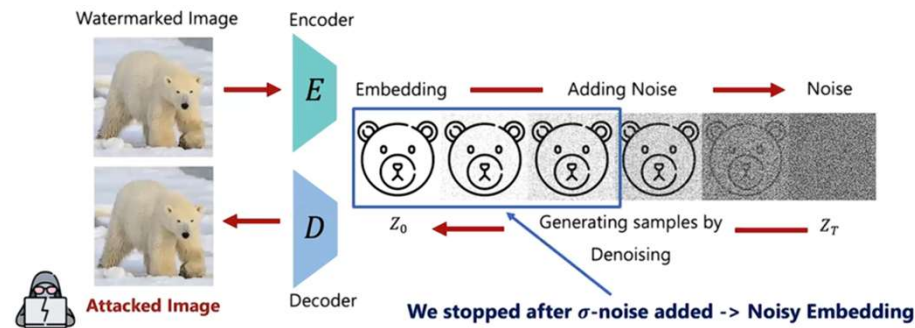
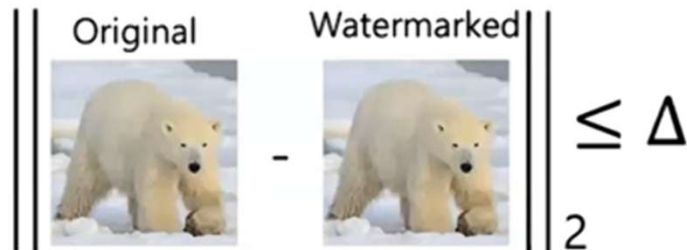
$$f(\varepsilon_1) = \Phi\left(\Phi^{-1}(1 - \varepsilon_1) - \frac{L_{x,w}\Delta}{\sigma}\right). \quad \varepsilon_2 \geq f(\varepsilon_1), \quad \forall \varepsilon_1 \in [0,1].$$



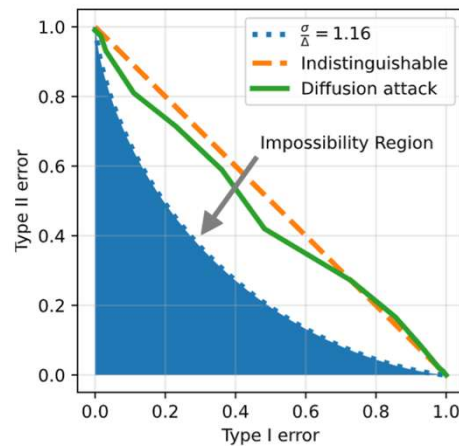
Background : Watermarking Schemes

$$\|x - x_w\| \leq \Delta.$$

$$|\phi(x_w) - \phi(x)| \leq L_{x,w}|x_w - x|$$



$$f(\varepsilon_1) = \Phi\left(\Phi^{-1}(1 - \varepsilon_1) - \frac{L_{x,w}\Delta}{\sigma}\right).$$

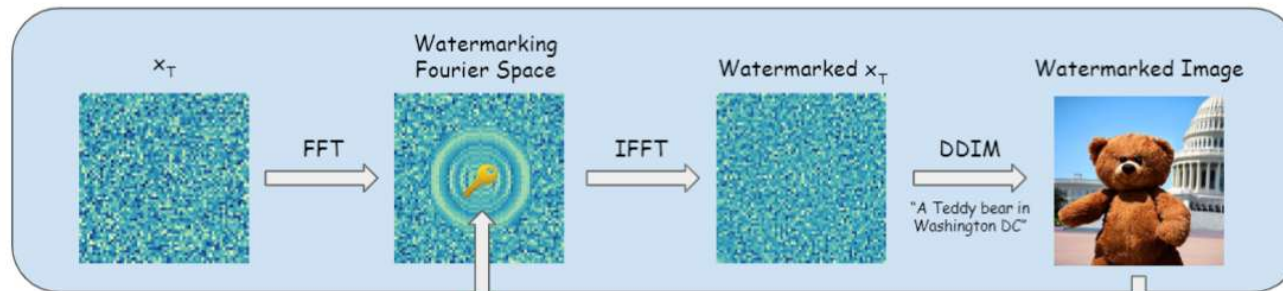


The regeneration attack is guaranteed to remove any pixel-based invisible watermarks.

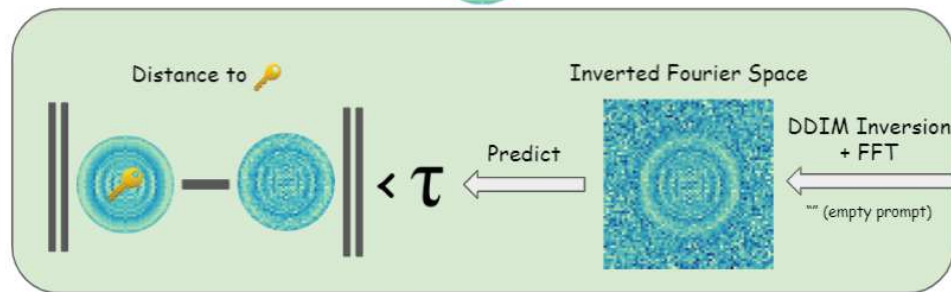
Background : Watermarking Schemes



Generation



Detection



Attack



Preliminaries

DDIM (deterministic strategy)

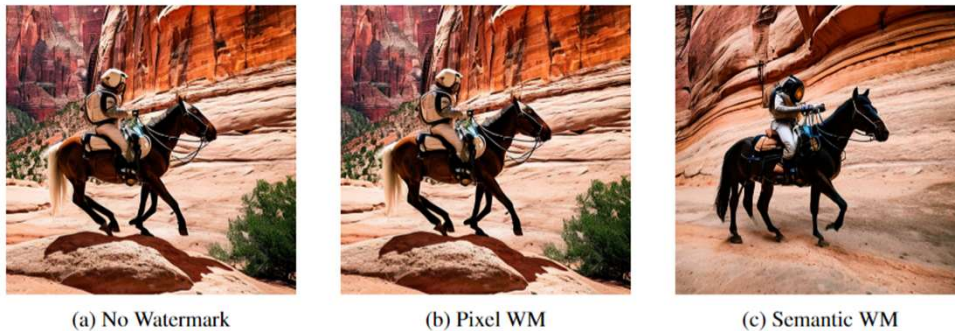
$$x_t = \sqrt{\alpha_t}x_0 + \sqrt{1 - \alpha_t}\epsilon$$

$$x_{t-1} - x_t \approx x_{t+1} - x_t$$

$$x_{t+1} = \sqrt{\alpha_{t+1}}\hat{x}_0 + \sqrt{1 - \alpha_{t+1}}\epsilon_\theta(x_t)$$

Background : Watermarking Schemes

“An astronaut riding a horse in Zion National Park.”



➔ The robustness comes at the price of more visible differences.

$\|x - x_w\| \leq \Delta$. ➔ **Not Satisfied!**

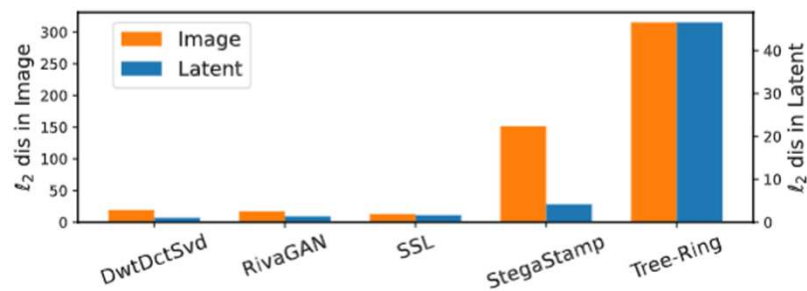
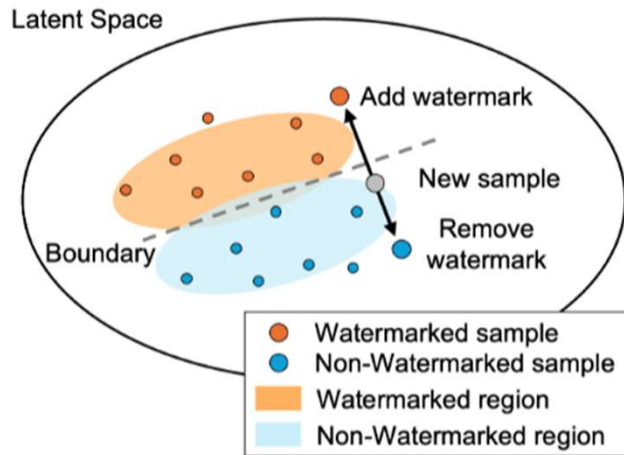
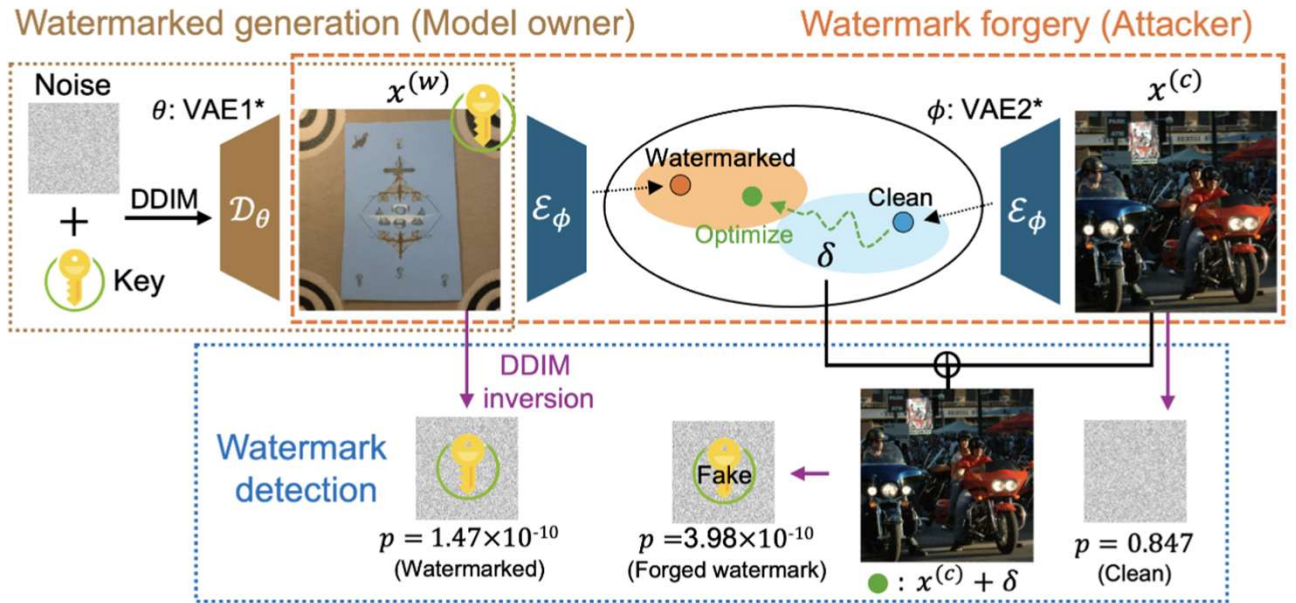
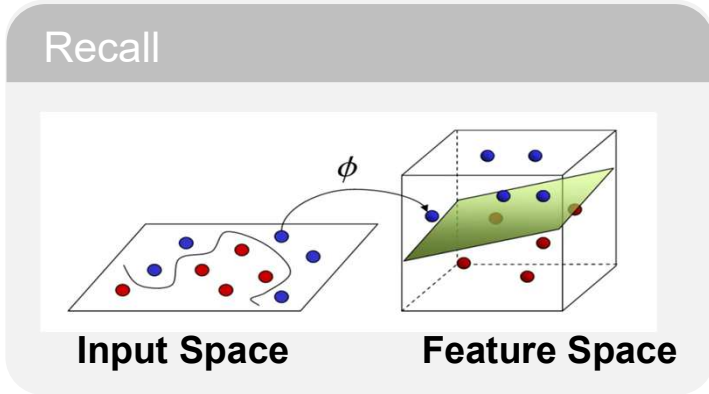


Table 2: Tree-Ring watermarks are robust against the regeneration attacks.

Attacker	MS-COCO TPR@FPR=0.01↓	SDP Generated TPR@FPR=0.01↓
Brightness-2	1.000	1.000
Contrast-2	1.000	1.000
JPEG-50	1.000	0.994
Gaussian noise-5	1.000	0.996
Gaussian blur-6	1.000	1.000
VAE-Bmshj2018-3	0.998	0.994
VAE-Cheng2020-3	1.000	0.994
Diffusion model-60	1.000	0.998

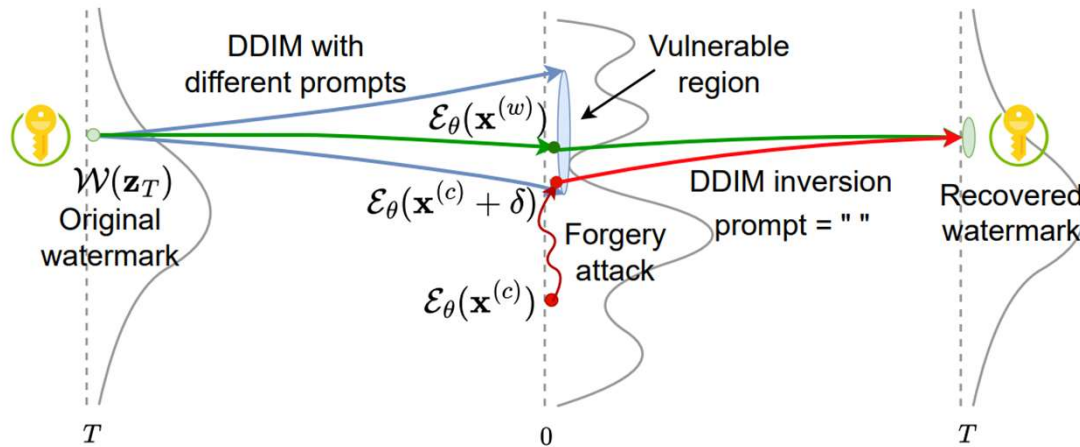
Motivation: Threat Model



* VAE1 and VAE2 can be different

Method: Watermark Region

Hypothesis: Latent Space Contains a “Watermark Region”



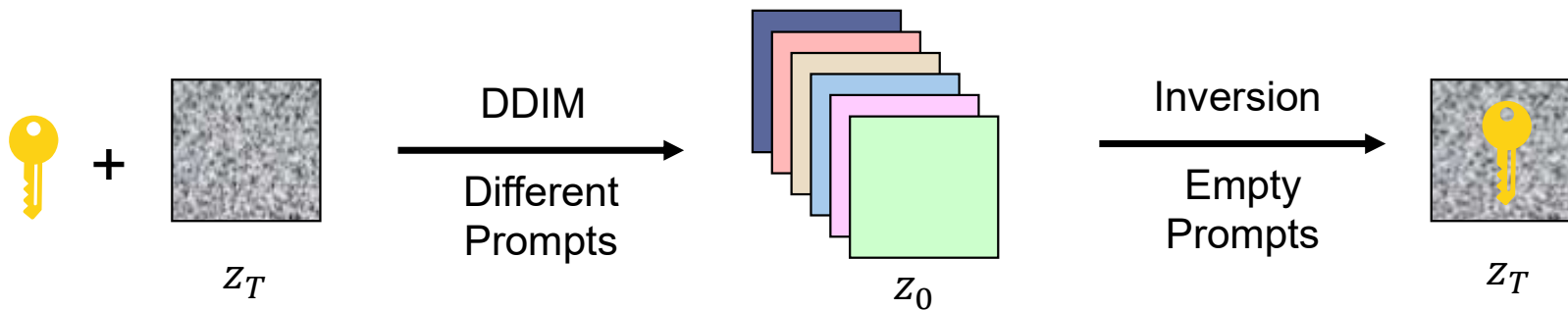
Preliminaries

DDIM (deterministic strategy)

$$x_t = \sqrt{\alpha_t}x_0 + \sqrt{1 - \alpha_t}\epsilon$$

$$x_{t-1} - x_t \approx x_{t+1} - x_t$$

$$x_{t+1} = \sqrt{\alpha_{t+1}}\hat{x}_0 + \sqrt{1 - \alpha_{t+1}}\epsilon_\theta(x_t)$$



Experiment: Watermark Region

Hypothesis: Latent Space Contains a “Watermark Region”

$$Z_0^{(w)}(\mathcal{W}, k) = \{z_0 \in Z_0 \mid \mathcal{M}_{\mathcal{W}}(\mathcal{J}^-(z_0), k) < \tau\}$$

Z_0 : Clean latent space

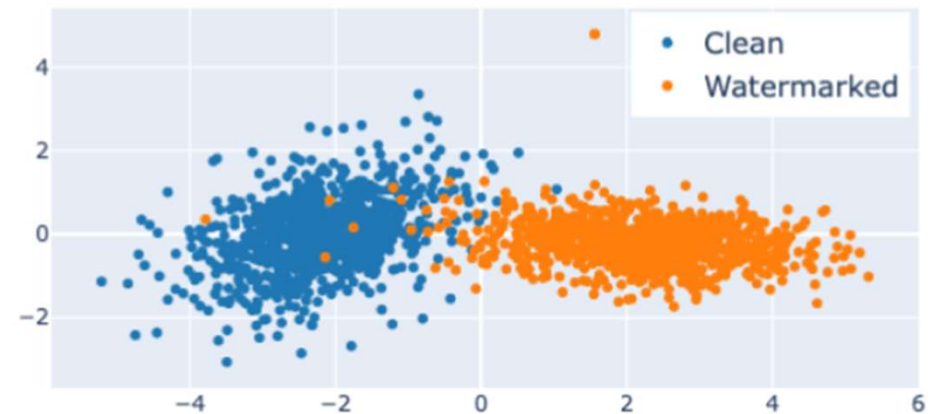
\mathcal{J}^- : DDIM inversion process

τ : Detection threshold (e.g., p-value < 0.05)

$\mathcal{M}_{\mathcal{W}}$: Matching function between the extracted key and the embedded key k

Experiment 1: Linear Separable

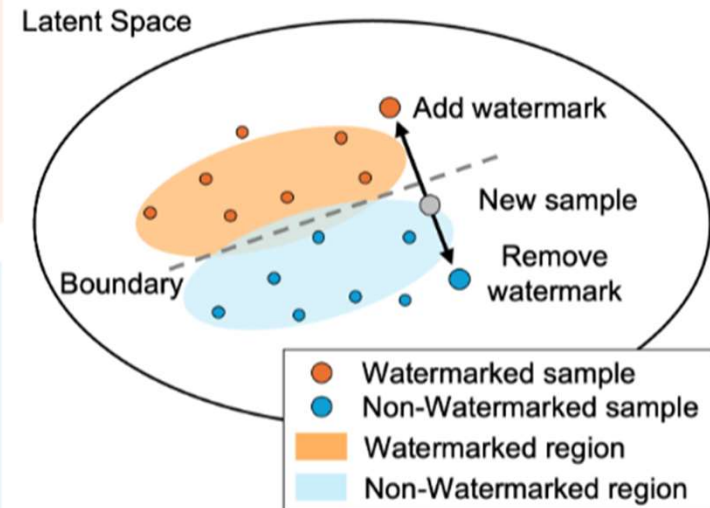
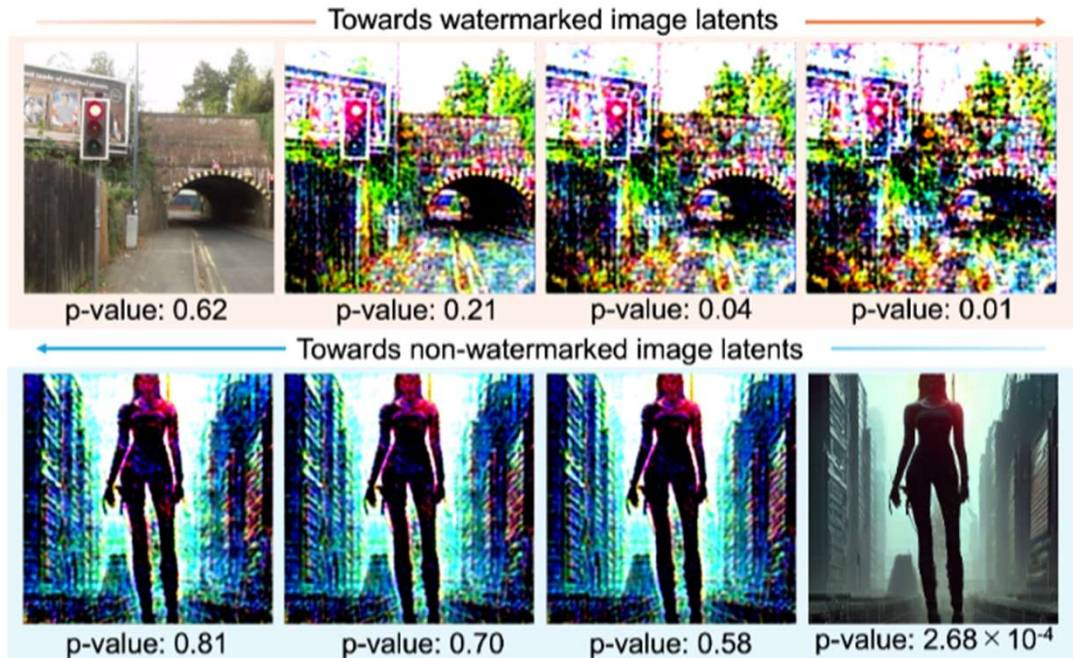
- 1000 watermarked images from a specific watermark and the same number of nonwatermarked images
- The linear SVM model yields an accuracy of 100% on the sampled images



Experiment : Watermark Region

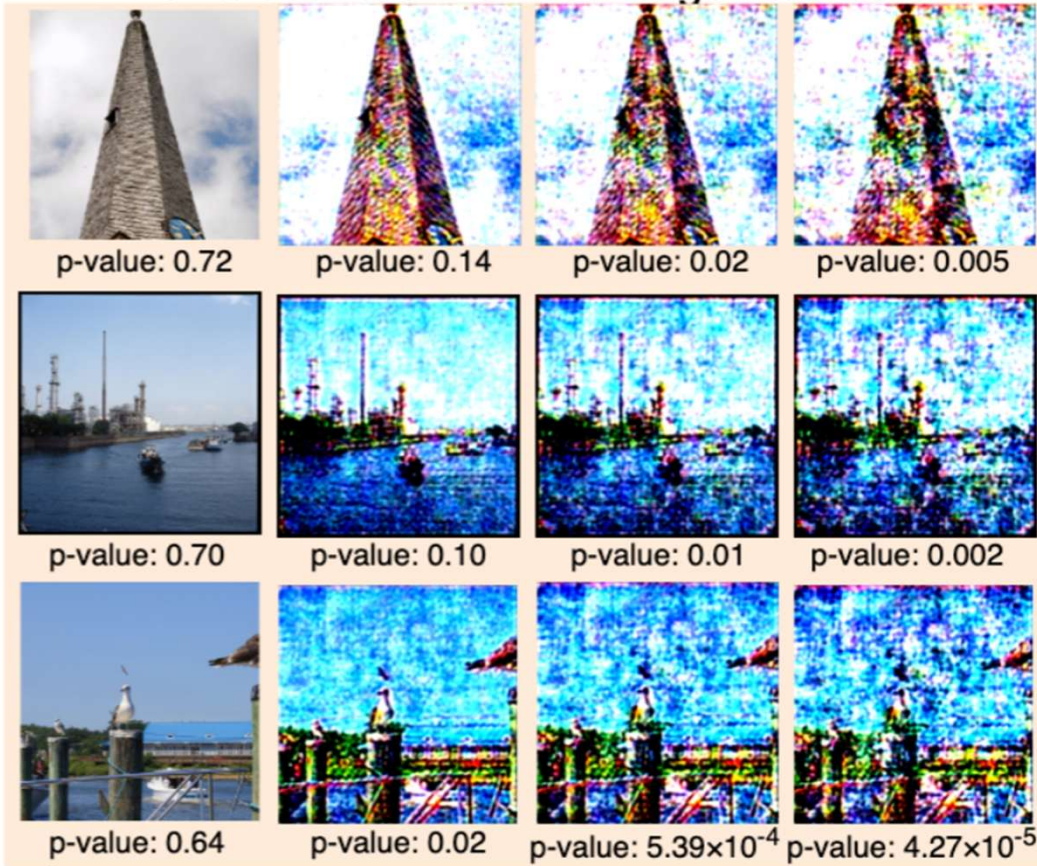
Hypothesis: Latent Space Contains a “Watermark Region”

Experiment 2: Coherent and Oriented latent subspace

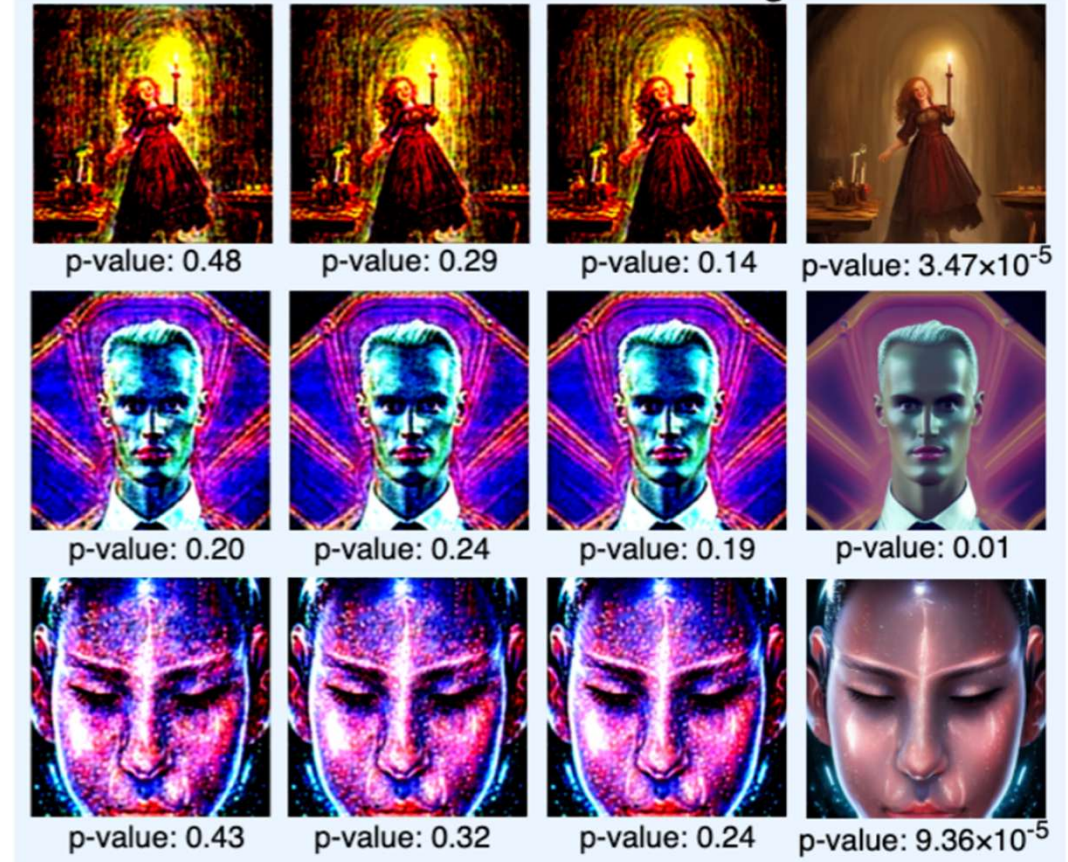


Experiment : Watermark Region

← Towards Watermarked Image Latents →



← Towards Non-Watermarked Image Latents →



Method: Imperceptible Attack

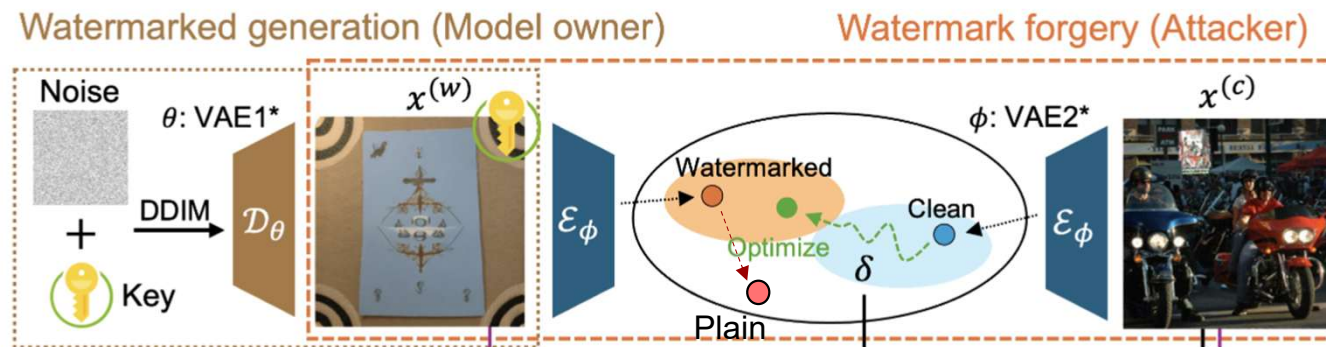
Forgery Attack

$$\min_{\delta} \|\mathcal{E}_{\phi}(x^{(c)} + \delta) - \mathcal{E}_{\phi}(x^{(w)})\|_2 + \lambda \|\delta\|_2$$

Removal Attack

$$\min_{\delta} \|\mathcal{E}_{\phi}(x^{(c)} + \delta) - \mathcal{E}_{\phi}(\mu_{x^{(w)}})\|_2 + \lambda \|\delta\|_2$$

$x^{(w)}$ watermarked image
 $x^{(c)}$ clean image
 $\mu_{x^{(w)}}$ plain image
 δ adversarial perturbation



Evaluation

Hypothesis

- The attacker has access to either
- (i) the VAE of the watermarked diffusion model
- (ii) a proxy VAE that was trained on a similar dataset

Diffusion Model

- Utilize the VAE from SDv1.4 to remove/forge the watermark generated from both SDv1.4 and SDv2.0

Watermarking Systems

- Tree-Ring, RingID, WIND, and Gaussian Shading

Evaluation

P-Value

- Null Hypothesis: not watermarked.
- The p-value represents the probability of observing the extracted key.
- A small p-value, we reject H_0 and conclude the image is watermarked.

Attack Success Rate (ASR)

- Forgery Attack: Percentage of p-value < 0.05 after attack.
- Removal Attack: Percentage of p-value ≥ 0.05 after attack.

Image Quality Metrics

- L2 distance, PSNR, SSIM, LPIPS.

Experiments

Trade-off between attack success rate (ASR) and imperceptibility.

Method	Model	λ	ASR	l_2	l_∞	LPIPS	SSIM	PSNR	FID
Tree-Ring [35]	SDv1.4	5×10^4	78.65	33.90	0.69	0.17	0.89	34.32	41.27
		2×10^4	86.93	48.42	0.89	0.26	0.82	31.20	61.18
		1×10^4	91.06	63.22	1.10	0.33	0.76	28.87	81.27
	SDv2.0	5×10^4	79.89	34.09	0.69	0.17	0.88	34.26	41.22
		2×10^4	90.72	48.83	0.91	0.26	0.82	31.11	61.75
		1×10^4	93.81	63.78	1.08	0.34	0.76	28.78	80.54
RingID [5]	SDv1.4	5×10^4	100.0	38.45	0.68	0.20	0.87	33.21	49.97
		2×10^4	100.0	55.20	0.86	0.30	0.80	30.06	75.55
		1×10^4	100.0	73.08	1.03	0.38	0.73	27.63	101.45
	SDv2.0	5×10^4	100.0	37.31	0.66	0.19	0.87	33.48	47.87
		2×10^4	100.0	53.94	0.84	0.29	0.80	30.27	72.08
		1×10^4	100.0	71.53	1.00	0.37	0.73	27.82	98.12
WIND [2]	SDv1.4	5×10^4	97.56	38.82	0.70	0.20	0.87	33.11	54.41
		2×10^4	97.56	56.13	0.89	0.29	0.80	29.88	82.68
		1×10^4	97.56	74.66	1.06	0.38	0.73	27.38	111.84
	SDv2.0	5×10^4	100.0	37.47	0.67	0.19	0.87	33.45	49.02
		2×10^4	100.0	54.18	0.84	0.28	0.80	30.23	74.69
		1×10^4	100.0	71.86	0.99	0.37	0.74	27.78	101.33
Gaussian Shading [39]	SDv1.4	5×10^4	96.85	37.27	0.70	0.19	0.87	33.48	46.93
		2×10^4	96.96	54.00	0.88	0.29	0.80	30.21	69.50
		1×10^4	96.96	71.97	1.05	0.37	0.73	27.64	93.85
	SDv2.0	5×10^4	100.0	36.78	0.66	0.19	0.87	33.60	46.20
		2×10^4	100.0	52.99	0.85	0.29	0.80	30.42	69.35
		1×10^4	100.0	69.83	1.02	0.37	0.74	28.02	93.04

Experiments



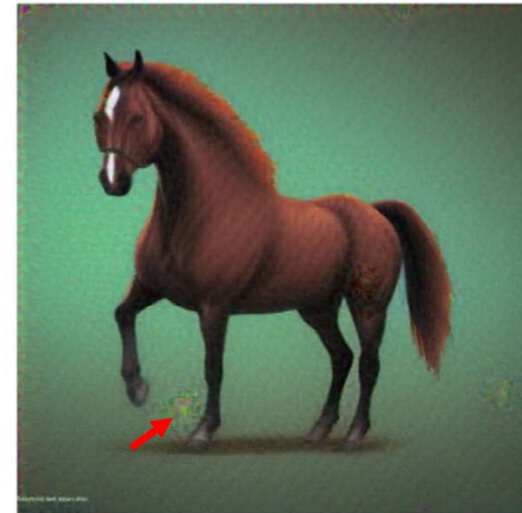
Original Watermarked



$\lambda = 5 \times 10^4$



$\lambda = 2 \times 10^4$



$\lambda = 1 \times 10^4$

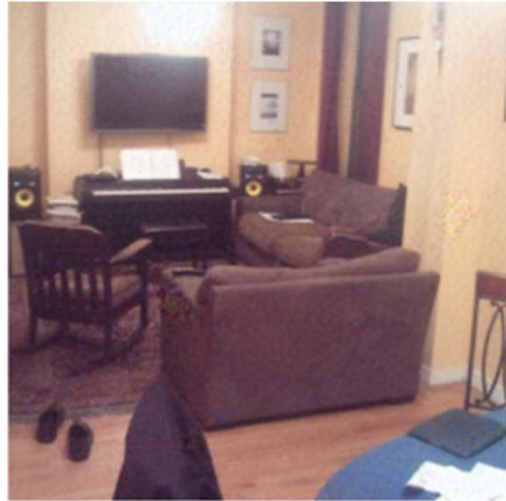
Tree-Ring watermarking method

λ	ASR	l_2	l_∞	LPIPS	SSIM	PSNR	FID
5×10^4	94.21	74.82	0.94	0.19	0.84	20.72	54.16
2×10^4	97.68	87.20	1.14	0.27	0.79	20.61	75.74
1×10^4	98.84	100.52	1.35	0.34	0.74	20.44	94.62

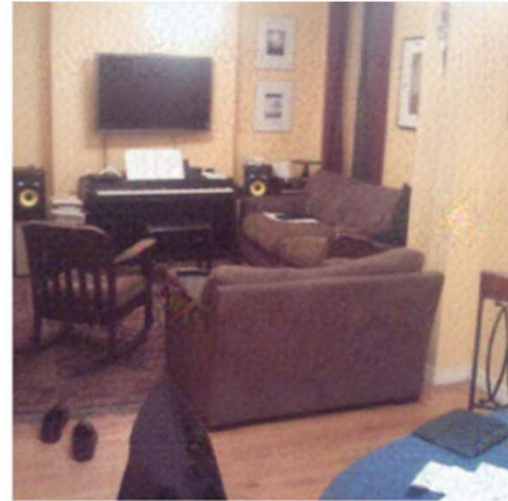
Experiments



Original



$\lambda = 5 \times 10^4$



$\lambda = 2 \times 10^4$



$\lambda = 1 \times 10^4$

Gaussian Shading
watermarking method

λ	ASR	l_2	l_∞	LPIPS	SSIM	PSNR	FID
5×10^4	96.85	37.27	0.70	0.19	0.87	33.48	46.93
2×10^4	96.96	54.00	0.88	0.29	0.80	30.21	69.50
1×10^4	96.96	71.97	1.05	0.37	0.73	27.64	93.85

Ablation

Different Optimization Objectives

Spatial Domain + L_2 Regularization:
$$\min_{\delta} \|\mathcal{E}_{\phi}(\mathbf{x}^{(c)} + \delta) - \mathcal{E}_{\phi}(\mathbf{x}^{(w)})\|_2 + \lambda \|\delta\|_2, \quad (3)$$

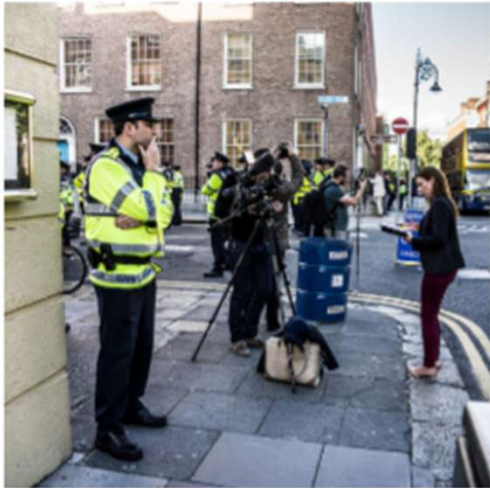
Spatial Domain + L_{∞} Constraint:
$$\min_{\delta} \|\mathcal{E}_{\phi}(\mathbf{x}^{(c)} + \delta) - \mathcal{E}_{\phi}(\mathbf{x}^{(w)})\|_2 \quad \text{s.t.} \quad \|\delta\|_{\infty} \leq \epsilon, \quad (5)$$

Frequency Domain + L_{∞} Constraint:
$$\min_{\delta} \|\mathcal{E}_{\phi}(\text{IDCT}(\text{DCT}(\mathbf{x}^{(c)}) + \mathbf{m} \odot \delta)) - \mathcal{E}_{\phi}(\mathbf{x}^{(w)})\|_2$$

$$\text{s.t.} \quad \|\delta\|_{\infty} \leq \epsilon. \quad (6)$$

Method	ASR	l_2	l_{∞}	LPIPS	SSIM	PSNR	FID
Eq. 5	83.24	68.04	0.10	0.35	0.66	28.31	69.35
Eq. 6	83.33	67.77	0.93	0.32	0.66	28.35	58.01
Eq. 3	84.91	44.34	0.85	0.23	0.84	31.97	56.30

Ablation



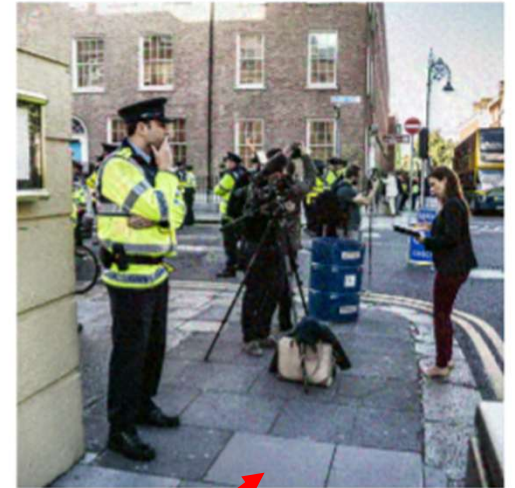
Original



Eq. 5



Eq. 6



Eq. 3

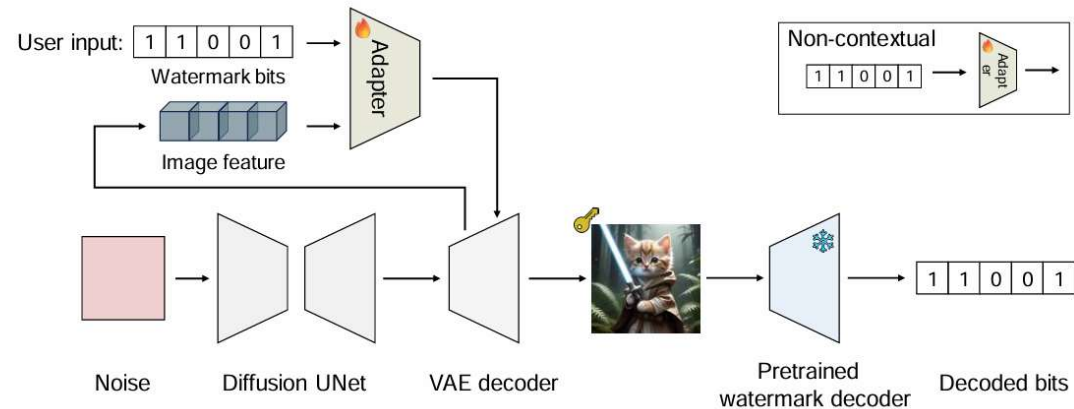
Limitations & Conclusion

- **Generative regeneration attacks:**

- Semantic Watermarks - Initial-Noise Watermarks

- **Latent-space adversarial attacks :**

- Rely on DDIM inversion watermarking methods
- Approaches that add a pattern during the latent decoding phase; require access to a similar decoder to attack the system



Limitations & Conclusion



- **Pixel-Level Watermarks:**
 - Perturbation limited to ℓ_2 ball → **generative regeneration**
- **Initial-Noise Watermarks:**
 - Rely on DDIM inversion → **latent-space adversarial attacks**



**Thanks for
Listening!**