

# The Stable Signature: Rooting Watermarks in Latent Diffusion Models

Pierre Fernandez<sup>1,2</sup> Guillaume Couairon<sup>1,3</sup> Hervé Jégou<sup>1</sup> Matthijs Douze<sup>1</sup> Teddy Furon<sup>2\*</sup>

<sup>1</sup>Meta AI

<sup>2</sup>Centre Inria de l’Université de Rennes

<sup>3</sup>Sorbonne University

## Abstract

Generative image modeling enables a wide range of applications but raises ethical concerns about responsible deployment. This paper introduces an active strategy combining image watermarking and Latent Diffusion Models. The goal is for all generated images to conceal an invisible watermark allowing for future detection and/or identification. The method quickly fine-tunes the latent decoder of the image generator, conditioned on a binary signature. A pre-trained watermark extractor recovers the hidden signature from any generated image and a statistical test then determines whether it comes from the generative model. We evaluate the invisibility and robustness of the watermarks on a variety of generation tasks, showing that Stable Signature works even after the images are modified. For instance, it detects the origin of an image generated from a text prompt, then cropped to keep 10% of the content, with 90+% accuracy at a false positive rate below  $10^{-6}$ .

## 1. Introduction

Recent progress in generative modeling and natural language processing enable easy creation and manipulation of photo-realistic images, such as with DALL-E 2 [61] or Stable Diffusion [65]. They have given birth to many image edition tools like ControlNet [101], Instruct-Pix2Pix [7], and others [13, 27, 68], that are establishing themselves as creative tools for artists, designers, and the general public.

While this is a great step forward for generative AI, it also renews concerns about undermining confidence in the authenticity or veracity of photo-realistic images. Indeed, methods to convincingly augment photo-realistic images have existed for a while, but generative AI significantly lowers the barriers to convincing synthetic image generation and edition (*e.g.* a generated picture recently won an art competition [28]). Not being able to identify that images are generated by AI makes it difficult to remove them from certain platforms and to ensure their compliance with ethical standards. It opens the doors to new risks like deep fakes, impersonation or copyright usurpation [8, 17].

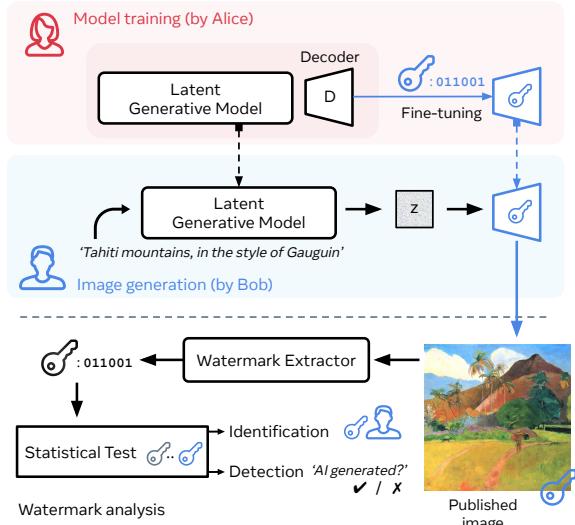


Figure 1. Overview. The latent decoder can be fine-tuned to preemptively embed a signature into all generated images.

A baseline solution to identify generated images is forensics, *i.e.* passive methods to detect generated/manipulated images. On the other hand, existing watermarking methods can be added on top of image generation. They are based on the idea of invisibly embedding a secret message into the image, which can then be extracted and used to identify the image. This has several drawbacks. If the model leaks or is open-sourced, the post-generation watermarking is easy to remove. The open source Stable Diffusion [67] is a case in point, since removing the watermark amounts to commenting out a single line in the source code.

Our Stable Signature method merges watermarking into the generation process itself, without any architectural changes. It adjusts the pre-trained generative model such that all the images it produces conceal a given watermark. There are several advantages to this approach [45, 96]. It protects both the generator and its productions. Besides, it does not require additional processing of the generated image, which makes the watermarking computationally lighter, straightforward, and secure. Model providers would then be able to deploy their models to different user groups with a unique watermark, and monitor that they are

\*Work supported by ANR / AID under Chaire SAIDA ANR-20-CHIA-0011. Correspondance to pfz@meta.com

used in a responsible manner. They could also give art platforms, news outlets and other sharing platforms the ability to detect when an image has been generated by their AI.

We focus on Latent Diffusion Models (LDM) [65] since they can perform a wide range of generative tasks. This work shows that simply fine-tuning a small part of the generative model – the decoder that generates images from the latent vectors – is enough to natively embed a watermark into all generated images. Stable Signature does not require any architectural change and does not modify the diffusion process. Hence it is compatible with most of the LDM-based generative methods [7, 13, 60, 68, 101]. The fine-tuning stage is performed by back-propagating a combination of a perceptual image loss and the hidden message loss from a watermark extractor back to the LDM decoder. We pre-train the extractor with a simplified version of the deep watermarking method HiDDeN [105].

We create an evaluation benchmark close to real world situations where images may be edited. The tasks are: detection of AI generated images, tracing models from their generations. For instance, we detect 90% of images generated with the generative model, even if they are cropped to 10% of their original size, while flagging only one false positive every  $10^6$  images. To ensure that the model’s utility is not weakened, we show that the FID [33] score of the generation is not affected and that the generated images are perceptually indistinguishable from the ones produced by the original model. This is done over several tasks involving LDM (text-to-image, inpainting, edition, etc.).

As a summary, (1) we efficiently merge watermarking into the generation process of LDMs, in a way that is compatible with most of the LDM-based generative methods; (2) we demonstrate how it can be used to detect and trace generated images, through a real-world evaluation benchmark; (3) we compare to post-hoc watermarking methods, showing that it is competitive while being more secure and efficient, and (4) evaluate robustness to intentional attacks.

## 2. Related Work

**Image generation** has long been dominated by GANs, still state-of-the-art on many datasets [36, 37, 38, 73, 84]. Transformers have also been successfully used for modeling image [63, 19] or video [76] distributions, providing higher diversity at the expense of increased inference time. Images are typically converted to token lists using vector-quantized architectures [20, 64, 93], relying on an image decoder.

Diffusion models [18, 35, 57, 77] have brought huge improvements in text-conditional image generation, being now able to synthesize high-resolution photo-realistic images for a wide variety of text prompts [3, 34, 62, 66, 71]. They can also perform conditional image generation tasks – like inpainting or text-guided image editing – by fine-tuning the diffusion model with additional conditioning,

e.g. masked input image, segmentation map, etc. [49, 70]. Because of their iterative denoising algorithm, diffusion models can also be adapted for image editing in a zero-shot fashion by guiding the generative process [13, 32, 39, 55, 82, 89]. All these methods, when applied on top of Stable Diffusion, operate in the latent space of images, requiring a latent decoder to produce an RGB image.

**Detection of AI-generated/manipulated images** is notably active in the context of deep-fakes [31, 104]. Many works focus on the detection of GAN-generated images [10, 30, 86, 103]. One way is to detect inconsistencies in the generated images, via lights, perspective or physical objects [21, 22, 44, 51, 85]. These approaches are restricted to photo-realistic images or faces but do not cover artworks where objects are not necessarily physically correct.

Other approaches use traces left by the generators in the spatial [53, 94] or frequency [25, 103] domains. They have extended to diffusion models in recent works [12, 74], and showed encouraging results. However purely relying on forensics and passive detection is limiting. As an example, the best performing method to our knowledge [12] is able to detect 50% of generated images for an FPR around 1/100. Put differently, if a user-generated content platform were to receive 1 billion images every day, it would need to wrongly flag 10 million images to detect only half of the generated images. Besides, passive techniques cannot trace images from different versions of the same model, conversely to active ones like watermarking.

**Image watermarking** has long been studied in the context of tracing and intellectual property protection [14]. More recently, deep learning encoder/extractor alternatives like HiDDeN [2, 43, 50, 97, 105] or iterative methods by Vukotić *et al.* [24, 41, 83] showed competitive results in terms of robustness to a wide range of transformations, namely geometric ones.

In the specific case of **generative models**, some works deal with watermarking the training set on which the generative model is learned [95]. It is highly inefficient since every new message to be embedded requires a new training pipeline. Merging the watermarking and the generative process is a recent idea [45, 90, 96, 99], that is closer to the model watermarking litterature [81]. They suffer from two strong limitations. First, these methods only apply to GAN, while LDM are beginning to replace them in most applications. Second, watermarking is incorporated in the training process of the GAN from the start. This strategy is very risky because the generative model training is more and more costly<sup>1</sup>. Our work shows that a quick fine-tuning of the latent decoder part of the generative model is enough to achieve a good watermarking performance, provided that the watermark extractor is well chosen.

<sup>1</sup>Stable Diffusion training costs ~\$600k of cloud compute ([Wikipedia](#)).

### 3. Problem Statement & Background

[Figure 1](#) shows a model provider *Alice* who deploys a latent diffusion model to users *Bobs*. Stable Signature embeds a binary signature into the generated images. This section derives how Alice can use this signature for two scenarios:

- *Detection*: “*Is it generated by my model?*”. Alice detects if an image was generated by her model. As many generations as possible should be flagged, while controlling the probability of flagging a natural image.
- *Identification*: “*Who generated this image?*”. Alice monitors who created each image, while avoiding to mistakenly identifying a Bob who did not generate the image.

#### 3.1. Image watermarking for detection

Alice embeds a  $k$ -bit binary signature into the generated images. The watermark extractor then decodes messages from the images it receives and detects when the message is close to Alice’s signature. An example application is to block AI-generated images on a content sharing platform.

**Statistical test** Let  $m \in \{0, 1\}^k$  be Alice’s signature. We extract the message  $m'$  from an image  $x$  and compare it to  $m$ . As done in previous works [45, 95], the detection test relies on the number of matching bits  $M(m, m')$ : if

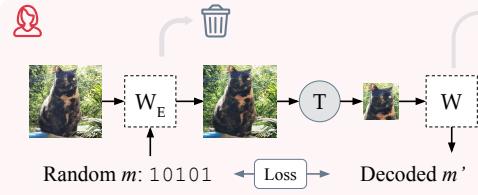
$$M(m, m') \geq \tau \text{ where } \tau \in \{0, \dots, k\}, \quad (1)$$

then the image is flagged. This provides a level of robustness to imperfections of the watermarking.

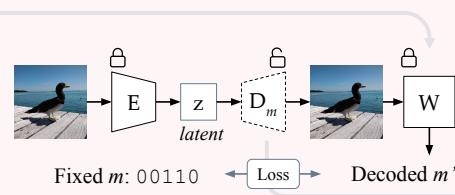
Formally, we test the statistical hypothesis  $H_1$ : “ $x$  was generated by Alice’s model” against the null hypothesis  $H_0$ : “ $x$  was not generated by Alice’s model”. Under  $H_0$  (*i.e.* for vanilla images), we assume that bits  $m'_1, \dots, m'_k$  are (*i.i.d.*) Bernoulli random variables with parameter 0.5. Then  $M(m, m')$  follows a binomial distribution with parameters  $(k, 0.5)$ . We verify this assumption experimentally in App. B.5. The False Positive Rate (FPR) is the probability that  $M(m, m')$  takes a value bigger than the threshold  $\tau$ . It is obtained from the CDF of the binomial distribution, and a closed-form can be written with the regularized incomplete beta function  $I_x(a; b)$ :

$$\text{FPR}(\tau) = \mathbb{P}(M > \tau | H_0) = I_{1/2}(\tau + 1, k - \tau). \quad (2)$$

(a) Pre-train watermark encoder/extractor



(b) Fine-tune LDM decoder



(c) Generate

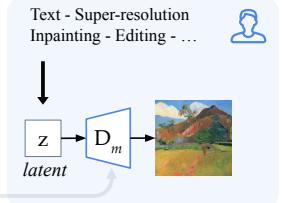


Figure 2. Steps of the method. (a) We pre-train a watermark encoder  $\mathcal{W}_E$  and extractor  $\mathcal{W}$ , to extract binary messages. (b) We fine-tune the decoder  $\mathcal{D}$  of the LDM’s auto-encoder with a fixed signature  $m$  such that all the generated images (c) lead to  $m$  through  $\mathcal{W}$ .

#### 3.2. Image watermarking for identification

Alice now embeds a signature  $m^{(i)}$  drawn randomly from  $\{0, 1\}^k$  into the model distributed to Bob<sup>(i)</sup> (for  $i = 1 \dots N$ , with  $N$  the number of Bobs). Alice can trace any misuse of her model: generated images violating her policy (gore content, deepfakes) are linked back to the specific Bob by comparing the extracted message to Bobs’ signatures.

**Statistical test** We compare the message  $m'$  from the watermark extractor to  $(m^{(1)}, \dots, m^{(N)})$ . There are now  $N$  detection hypotheses to test. If the  $N$  hypotheses are rejected, we conclude that the image was not generated by any of the models. Otherwise, we attribute the image to  $\text{argmax}_{i=1..N} M(m', m^{(i)})$ . With regards to the detection task, false positives are more likely since there are  $N$  tests. The global FPR at a given threshold  $\tau$  is:

$$\text{FPR}(\tau, N) = 1 - (1 - \text{FPR}(\tau))^N \approx N \cdot \text{FPR}(\tau). \quad (3)$$

Equation (3) (resp. (2)), is used reversely: we find threshold  $\tau$  to achieve a required FPR for identification (resp. detection). Note that these formulae hold only under the assumption of i.i.d. Bernoulli bits extracted from vanilla images. This crucial point is enforced in the next section.

### 4. Method

Stable Signature modifies the generative network so that the generated images have a given signature through a fixed watermark extractor. It is trained in two phases. First, we create the watermark extractor network  $\mathcal{W}$ . Then, we fine-tune the Latent Diffusion Model (LDM) decoder  $\mathcal{D}$ , such that all generated images have a given signature through  $\mathcal{W}$ .

#### 4.1. Pre-training the watermark extractor

We use HiDDen [105], a classical method in the deep watermarking literature. It jointly optimizes the parameters of watermark encoder  $\mathcal{W}_E$  and extractor network  $\mathcal{W}$  to embed  $k$ -bit messages into images, robustly to transformations that are applied during training. We discard  $\mathcal{W}_E$  after training, since only  $\mathcal{W}$  serves our purpose.

Formally,  $\mathcal{W}_E$  takes as inputs a cover image  $x_o \in \mathbb{R}^{W \times H \times 3}$  and a  $k$ -bit message  $m \in \{0, 1\}^k$ . Similar to

ReDMark [2],  $\mathcal{W}_E$  outputs a residual image  $\delta$  of the same size as  $x_o$ , that is multiplied by a factor  $\alpha$  to produce watermarked image  $x_w = x_o + \alpha\delta$ . At each optimization step an image transformation  $T$  is sampled from a set  $\mathcal{T}$  that includes common image processing operations such as cropping and JPEG compression<sup>2</sup>. A “soft” message is extracted from the transformed image:  $m' = \mathcal{W}(T(x_w))$  (at inference time, the decoded message is given by the signs of the components of  $m'$ ). The *message loss* is the Binary Cross Entropy (BCE) between  $m$  and the sigmoid  $\sigma(m')$ :

$$\mathcal{L}_m = - \sum_{i=1}^k m_i \cdot \log \sigma(m'_i) + (1 - m_i) \cdot \log(1 - \sigma(m'_i)).$$

The network architectures are kept simple to ease the LDM fine-tuning in the second phase. They are the same as HiDDeN [105] (see App. A.1) with two changes.

First, since  $\mathcal{W}_E$  is discarded, its perceptual quality is not as important, so the perceptual loss or the adversarial network are not needed. Instead, the distortion is constrained by a tanh function on output of  $\mathcal{W}_E$  and by the scaling factor  $\alpha$ . This improves the bit accuracy of the recovered message and makes it possible to increase its size  $k$ .

Second, we observed that  $\mathcal{W}$ ’s output bits for vanilla images are correlated and highly biased, which violates the assumptions of Sec. 3.1. Therefore we remove the bias and decorrelate the outputs of  $\mathcal{W}$  by applying a PCA whitening transformation (more details in App. A.1).

## 4.2. Fine-tuning the generative model

In LDM, the diffusion happens in the latent space of an auto-encoder. The latent vector  $z$  obtained at the end of the diffusion is input to decoder  $\mathcal{D}$  to produce an image. Here we fine-tune  $\mathcal{D}$  such that the image contains a given message  $m$  that can be extracted by  $\mathcal{W}$ . Stable Signature is compatible with many generative tasks, since modifying only  $\mathcal{D}$  does not affect the diffusion process.

First, we fix the signature  $m = (m_1, \dots, m_k) \in \{0, 1\}^k$ . The fine-tuning of  $\mathcal{D}$  into  $\mathcal{D}_m$  is inspired by the original training of the auto-encoder in LDM [65].

Training image  $x \in \mathbb{R}^{H \times W \times 3}$  is fed to the LDM encoder  $\mathcal{E}$  that outputs activation map  $z = \mathcal{E}(x) \in \mathbb{R}^{h \times w \times c}$ , downsampled by a power-of-two factor  $f = H/h = W/w$ . The decoder reconstructs an image  $x' = \mathcal{D}_m(z)$  and the extractor recovers  $m' = \mathcal{W}(x')$ . The *message loss* is the BCE between  $m'$  and the original  $m$ :  $\mathcal{L}_m = \text{BCE}(\sigma(m'), m)$ .

In addition, the original decoder  $\mathcal{D}$  reconstructs the image without watermark:  $x'_o = \mathcal{D}(z)$ . The *image perceptual loss*  $\mathcal{L}_i$  between  $x'$  and  $x'_o$ , controls the distortion. We use the Watson-VGG perceptual loss introduced by Czolbe *et al.* [15], an improved version of LPIPS [102]. It is essential

<sup>2</sup>The transformation needs to be differentiable in pixel space. This is not the case for JPEG compression so we use the forward attack simulation layer introduced by Zhang *et al.* [98].

that the decoder learns luminance and contrast masking to add less perceivable watermarks.

The weights of  $\mathcal{D}_m$  are optimized in a few backpropagation steps to minimize

$$\mathcal{L} = \mathcal{L}_m + \lambda_i \mathcal{L}_i. \quad (4)$$

This is done over 100 iterations with the AdamW optimizer [48] and batch of size 4, *i.e.* the fine-tuning sees *less than 500 images* and takes *one minute on a single GPU*. The learning rate follows a cosine annealing schedule with 20 iterations of linear warmup to  $10^{-4}$  and decays to  $10^{-6}$ . The factor  $\lambda_i$  in (4) is set to 2.0 by default.

## 5. Text-to-Image Watermarking Performance

This section shows the potential of our method for detection and identification of images generated by a Stable-Diffusion-like model [65]<sup>3</sup>. We apply generative models watermarked with 48-bit signatures on prompts of the MS-COCO [46] validation set. We evaluate detection and identification on the outputs, as illustrated in Figure 1.

We evaluate their robustness to different transformations applied to generated images: strong cropping (10% of the image remaining), brightness shift (strength factor 2.0), as well as a combination of crop 50%, brightness shift 1.5 and JPEG 80. This covers typical geometric and photometric edits (see Fig. 5 for visual examples).

The performance is partly obtained from experiments and partly by extrapolating small-scale measurements.

### 5.1. Detection results

For detection, we fine-tune the decoder of the LDM with a random key  $m$ , generate 1000 images and use the test of Eq. (1). We report the tradeoff between True Positive Rate (TPR), *i.e.* the probability of flagging a generated image and the FPR, while varying  $\tau \in \{0, \dots, 48\}$ . For instance, for  $\tau = 0$ , we flag all images so  $\text{FPR} = 1$ , and  $\text{TPR} = 1$ . The TPR is measured directly. In contrast the FPR is inferred from Eq. (2), because it would otherwise be too small to be measured on reasonably sized problems (this approximation is validated experimentally in App. B.6). The experiment is run on 10 random signatures and we report averaged results.

Figure 3 shows the tradeoff under image transformations. For example, when the generated images are not modified, Stable Signature detects 99% of them, while only 1 vanilla image out of  $10^9$  is flagged. At the same  $\text{FPR} = 10^{-9}$ , Stable Signature detects 84% of generated images for a crop that keeps 10% of the image, and 65% for a transformation that combines a crop, a color shift, and a JPEG compression. For comparison, we report results of a state-of-the-art passive method [12], applied on resized

<sup>3</sup>We refrain from experimenting with pre-existing third-party generative models, such as Stable Diffusion or LDMs, and instead use a large diffusion model (2.2B parameters) trained on an internal dataset of 330M licensed image-text pairs.

and compressed images. As to be expected, we observe that these baseline results have orders of magnitudes larger FPR than Stable Signature, which actively marks the content.

## 5.2. Identification results

Each Bob has its own copy of the generative model. Given an image, the goal is to find if any of the  $N$  Bobs created it (detection) and if so, which one (identification). There are 3 types of error: *false positive*: flag a vanilla image; *false negative*: miss a generated image; *false accusation*: flag a generated image but identify the wrong user.

For evaluation, we fine-tune  $N' = 1000$  models with random signatures. Each model generates 100 images. For each of these 100k watermarked images, we extract the Stable Signature message, compute the matching score with all  $N$  signatures and select the user with the highest score. The image is predicted to be generated by that user if this score is above threshold  $\tau$ . We determined  $\tau$  such that  $\text{FPR} = 10^{-6}$ , see Eq. (3). For example, for  $N = 1$ ,  $\tau = 41$  and for  $N = 1000$ ,  $\tau = 44$ . Accuracy is extrapolated beyond the  $N'$  users by adding additional signatures and having  $N > N'$  (*e.g.* users that have not generated any images).

**Figure 4** reports the per-transformation identification accuracy. For example, we identify a user among  $N=10^5$  with 98% accuracy when the image is not modified. Note that for the combined edit, this becomes 40%. This may still be dissuasive: if a user generates 3 images, he will be identified 80% of the time. We observe that at this scale, the false accusation rate is zero, *i.e.* we never identify the wrong user. This is because  $\tau$  is set high to avoid FPs, which also makes false accusations unlikely. We observe that the identification accuracy decreases when  $N$  increases, because the threshold  $\tau$  required to avoid false positives is higher when  $N$  increases, as pointed out by the approximation in (3). In a nutshell, by distributing more models, Alice trades some accuracy of detection against the ability to identify users.

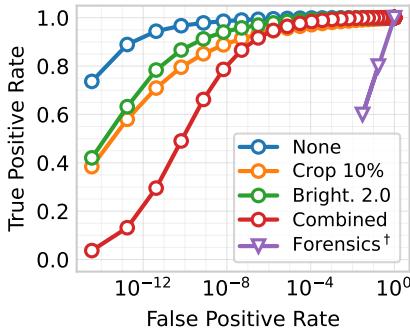


Figure 3. **Detection results.** TPR/FPR curve of the detection under different transformations.  $\text{Forensics}^\dagger$  indicates passive detection (without watermark) [12].

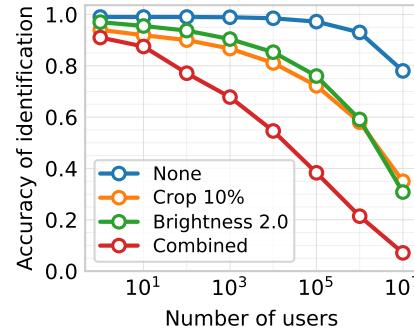


Figure 4. **Identification results.** Proportion of well-identified users. Detection with  $\text{FPR}=10^{-6}$  is run beforehand, and we consider it an error if the image is not flagged.

## 6. Experimental Results

We presented in the previous section how to leverage watermarks for detection and identification of images generated from text prompts. We now present more general results on robustness and image quality for different generative tasks. We also compare Stable Signature to other watermarking algorithms applied post-generation.

### 6.1. Tasks & evaluation metrics

Since our method only involves the LDM decoder, it makes it compatible with many generative tasks. We evaluate text-to-image generation and image edition on the validation set of MS-COCO [46], super-resolution and inpainting on the validation set of ImageNet [16] (all evaluation details are available in App. A.3).

We evaluate the image distortion with the Peak Signal-to-Noise Ratio (PSNR), which is defined as  $\text{PSNR}(x, x') = -10 \cdot \log_{10}(\text{MSE}(x, x'))$ , for  $x, x' \in [0, 1]^{c \times h \times w}$ , as well as Structural Similarity score (SSIM) [87]. They compare images generated with and without watermark. On the other hand, we evaluate the diversity and quality of the generated images with the Fréchet Inception Distance (FID) [33]. The bit accuracy – the percentage of bits correctly decoded – evaluates the watermarks’ robustness.

### 6.2. Image generation quality

**Figure 6** shows qualitative examples of how the image generation is altered by the latent decoder’s fine-tuning. The difference is very hard to perceive even for a trained eye. This is surprising for such a low PSNR, especially since the watermark embedding is not constrained by any Human Visual System like in professional watermarking techniques. Most interestingly, the LDM decoder has indeed learned to add the watermark signal only over textured areas where the human eyes are not sensitive, while the uniform backgrounds are kept intact (see the pixel-wise difference).



Figure 5. Transformations evaluated in Sec. 5 & 6. ‘Combined’ is made of crop 50%, brightness adjustment 1.5 and JPEG 80 compression.

Tasks		LDM [65]	PSNR / SSIM ↑		FID ↓	Bit accuracy ↑ on:			
			None	Crop	Brigh.	Comb.			
Text-to-Image	LDM [65]	30.0 / 0.89	19.6 ( <b>-0.3</b> )	0.99	0.95	0.97	0.92		
	DiffEdit [13]	31.2 / 0.92	15.0 ( <b>-0.3</b> )	0.99	0.95	0.98	0.94		
	Glide [56]	31.1 / 0.91	16.8 ( <b>+0.6</b> )	0.99	0.97	0.98	0.93		
		37.8 / 0.98	9.0 ( <b>+0.1</b> )	0.89	0.76	0.84	0.78		
WM Methods	Super-Resolution	LDM [65]	34.0 / 0.94	11.6 ( <b>+0.0</b> )	0.98	0.93	0.96	0.92	
	<i>Post generation</i>								
	Dct-Dwt [14]	0.14 (s/img)	39.5 / 0.97	19.5 ( <b>-0.4</b> )	0.86	0.52	0.51	0.51	
	SSL Watermark [24]	0.45 (s/img)	31.1 / 0.86	20.6 ( <b>+0.7</b> )	1.00	0.73	0.93	0.66	
<i>Merged in generation</i>	FNNS [41]	0.28 (s/img)	32.1 / 0.90	19.0 ( <b>-0.9</b> )	0.93	0.93	0.91	0.93	
	HiDDeN [105]	0.11 (s/img)	32.0 / 0.88	19.7 ( <b>-0.2</b> )	0.99	0.97	0.99	0.98	
	Stable Signature	0.00 (s/img)	30.0 / 0.89	19.6 ( <b>-0.3</b> )	0.99	0.95	0.97	0.92	

**Table 1** presents a quantitative evaluation of image generation quality on the different tasks. We report the FID, and the average PSNR and SSIM that are computed between the images generated by the fine-tuned LDM and the original one. The results show that no matter the task, the watermarking has very small impact on the FID of the generation.

The average PSNR is around 30 dB and SSIM around 0.9 between images generated by the original and a watermarked model. They are a bit low from a watermarking perspective because we do not explicitly optimize for them. Indeed, in a real world scenario, one would only have the watermarked version of the image. Therefore we don’t need to be as close as possible to the original image but only want to generate artifacts-free images. Without access to the image generated by the original LDM, it is very hard to tell whether a watermark is present or not.

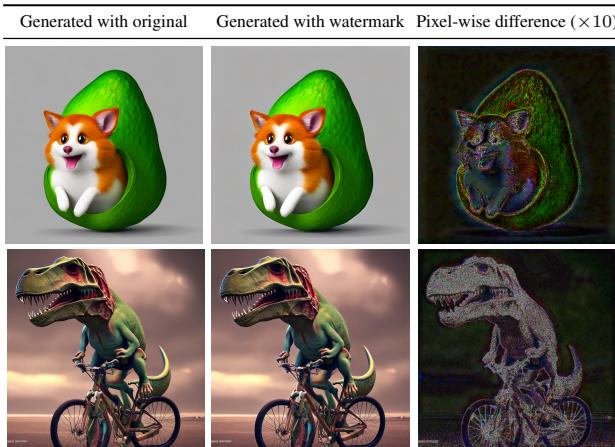


Figure 6. Images generated with Stable Diffusion. The PSNR is 35.4 dB in the first row and 28.6 dB in the second. Images generated with Stable Signature look natural because modified areas are located where the eye is not sensitive. More examples in App. C.

**Table 1.** Generation quality and comparison to post-hoc watermarking on  $512 \times 512$  images and 48-bit signatures. PSNR and SSIM are computed between generations of the original and watermarked generators. For FID, we show in (color) the difference with regards to original. Post-hoc watermarking is evaluated on text-generated images. (App. B.2 gives results on more transformations, and App. A gives more details on the evaluations.) Overall, Stable Signature has minimal impact on generation quality. It has comparable robustness to post-hoc methods while being rooted in the generation itself.

**Table 2.** Watermark robustness on image transformations applied before decoding, details of which are available in App. A.3. We report the bit accuracy, averaged over  $10 \times 1k$  images generated from COCO prompts with 10 different keys.

Attack	Bit acc.	Comb.	0.92	Sharpness 2.0	0.99	
None	0.99	Bright.	2.0	0.97	Med. Filter k=7	0.94
Crop 0.1	0.95	Cont.	2.0	0.98	Resize 0.7	0.91
JPEG 50	0.88	Sat.	2.0	0.99	Text overlay	0.99

### 6.3. Watermark robustness

We evaluate the robustness of the watermark to different image transformations applied before extraction. For each task, we generate 1k images with 10 models fine-tuned for different messages, and report the average bit accuracy in **Table 1**. Additionally, **Table 2** reports results on more image transformations for images generated from COCO prompts. The main evaluated transformations are presented in Fig. 5 (more evaluation details are available in App. A.3).

We see that the watermark is indeed robust for several tasks and across transformations. The bit accuracy is always above 0.9, except for inpainting, when replacing only the masked region of the image (between 1–50% of the image, with an average of 27% across masks). Besides, the bit accuracy is not perfect even without edition, mainly because there are images that are harder to watermark (*e.g.* the ones that are very uniform, like the background in Fig. 6) and for which the accuracy is lower.

Note that the robustness comes even without any transformation during the LDM fine-tuning phase: it is due to the watermark extractor. If the watermark embedding pipeline is learned to be robust against an augmentation, then the LDM will learn how to produce watermarks that are robust against it during fine-tuning.

### 6.4. Comparison to post-hoc watermarking

An alternative way to watermark generated images is to process them after the generation (post-hoc). This may

Table 3. Role of the attack simulation layer.

$\lambda_i$ for fine-tuning	0.8	0.4	0.2	0.1	0.05	0.025
PSNR $\uparrow$	31.4	30.6	29.7	28.5	26.8	24.6
Bit acc. $\uparrow$ on ‘comb.’	0.85	0.88	0.90	0.92	0.94	0.95

Table 4. Role of the attack simulation layer.

Seen at $\mathcal{W}$ training	Bit accuracy $\uparrow$ at test time:				
	Crop 0.1	Rot. 90	JPEG 50	Bright. 2.0	Res. 0.7
✗	1.00	0.56	0.50	0.99	0.48
✓	1.00	0.99	0.90	0.99	0.91

be simpler, but less secure and efficient than Stable Signature. We compare our method to a frequency based method, DCT-DWT [14], iterative approaches (SSL Watermark [24] and FNNS [41]), and an encoder/decoder one like HiD-DeN [105]. We choose DCT-DWT since it is employed by the original open source release of Stable Diffusion [67], and the other methods because of their performance and their ability to handle arbitrary image sizes and number of bits. We use our implementations (see details in App. A.4).

Table 1 compares the generation quality and the robustness over 5k generated images. Overall, Stable Signature achieves comparable results in terms of robustness. HiD-DeN’s performance is a bit higher but its output bits are not i.i.d. meaning that it cannot be used with the same guarantees as the other methods. We also observe that post-hoc generation gives worse qualitative results, images tend to present artifacts (see Fig. 13 in the supplement). One explanation is that Stable Signature is merged into the high-quality generation process with the LDM auto-encoder model, which is able to modify images in a more subtle way.

## 6.5. Can we trade image quality for robustness?

We can choose to maximize the image quality or the robustness of the watermark thanks to the weight  $\lambda_i$  of the perceptual loss in (4). We report the average PSNR of 1k generated images, as well as the bit accuracy obtained on the extracted message for the ‘Combined’ editing applied before detection (qualitative results are in App. B.1). A higher  $\lambda_i$  leads to an image closer to the original one, but to lower bit accuracies on the extracted message, see Table 3.

## 6.6. Attack simulation layer

Watermark robustness against image transformations depends solely on the watermark extractor. here, we pre-train them with or without specific transformations in the simulation layer, on a shorter schedule of 50 epochs, with  $128 \times 128$  images and 16-bits messages. From there, we plug them in the LDM fine-tuning stage and we generate 1k images from text prompts. We report the bit accuracy of the extracted watermarks in Table 4. The extractor is naturally robust to some transformations, such as crops or brightness, without being trained with them, while others, like rotations

or JPEG, require simulation during training for the watermark to be recovered at test time. Empirically we observed that adding a transformation improves results for the latter, but makes training more challenging.

## 7. Attacks on Stable Signature’s Watermarks

We examine the watermark’s resistance to intentional tampering, as opposed to distortions that happen without bad intentions like crops or compression (discussed in Sec. 5). We consider two threat models: one is typical for many image watermarking methods [14] and operates at the image level, and another targets the generative model level. For image-level attacks, we evaluate on 5k images generated from COCO prompts. Full details on the following experiments can be found in Appendix A.5.

### 7.1. Image-level attacks

**Watermark removal.** Bob alters the image to remove the watermark with deep learning techniques, like methods used for adversarial purification [75, 91] or neural auto-encoders [1, 47]. Note that this kind of attacks has not been explored in the image watermarking literature to our knowledge. Figure 7 evaluates the robustness of the watermark against neural auto-encoders [4, 11, 20, 65] at different compression rates. To reduce the bit accuracy closer to random (50%), the image distortion needs to be strong (PSNR<26). However, assuming the attack is *informed on the generative model*, i.e. the auto-encoder is the same as the one used to generate the images, the attack becomes much more effective. It erases the watermark while achieving high quality (PSNR>29). This is because the image is modified precisely in the bandwidth where the watermark

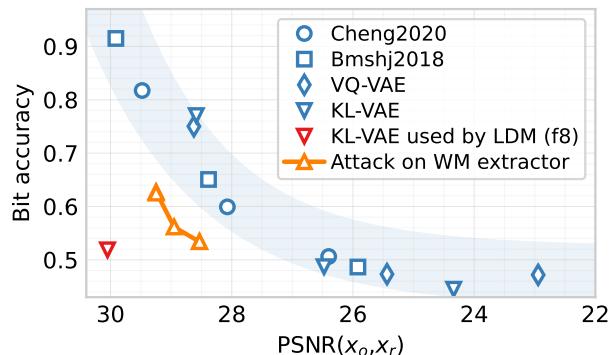


Figure 7. **Removal attacks.**  $x_o$  is the image produced by the original generator,  $x_r$  is the version produced by the watermarked generator and then attacked. Bit accuracy is on the watermark extracted from  $x_r$ . Neural auto-encoders follow the same trend, except with the one used by the LDM (‘KL-f8’ for our LDM). When access to the watermark extractor is granted, adversarial attacks also remove the watermark at lower PSNR budget.

is embedded. Note that this assumption is strong, because Alice does not need to distribute the original generator.

**Watermark removal & embedding (white-box).** To go further, we assume that the attack is *informed on the watermark extractor* – e.g. because it has leaked. Bob can use an adversarial attack to remove the watermark by optimizing the image under a PSNR constraint. The objective is to minimize the  $\ell_2$  distance between a random binary message sampled beforehand and the extractor’s output, effectively replacing the original signature with a random one. It makes it possible to erase the watermark with a lower distortion budget, as seen in Fig. 7.

Instead of removing the watermark, an attacker could embed a signature into vanilla images (unauthorized embedding [14]) to impersonate another Bob of whom they have a generated image. It highlights the importance of keeping the watermark extractor private.

## 7.2. Network-level attacks

**Model purification.** Bob gets Alice’s generative model and uses a fine-tuning process akin to Sec. 4.2 to eliminate the watermark embedding – that we coin *model purification*. This involves removing the message loss  $\mathcal{L}_m$ , and shifting the focus to the perceptual loss  $\mathcal{L}_i$  between the original image and the one reconstructed by the LDM auto-encoder.

Figure 8 shows the results of this attack for the MSE loss. The PSNR between the watermarked and purified images is plotted at various stages of fine-tuning. Empirically, it is difficult to significantly reduce the bit accuracy without compromising the image quality: artifacts start to appear during the purification.

**Model collusion.** Users may collude by aggregating their models. For instance,  $\text{Bob}^{(i)}$  and  $\text{Bob}^{(j)}$  can average the weights of their models (like Model soups [88]) creating a new model to deceive identification. We found that the bit at position  $\ell$  output by the extractor will be 0 (resp. 1) when the  $\ell$ -th bits of  $\text{Bob}^{(i)}$  and  $\text{Bob}^{(j)}$  are both 0 (resp. 1), and that the extracted bit is random when their bits disagree. We

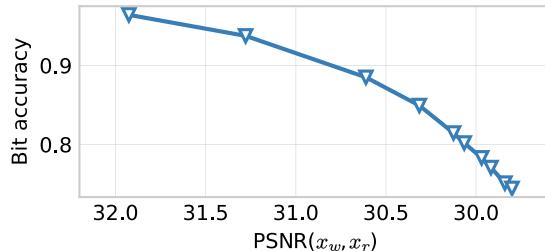
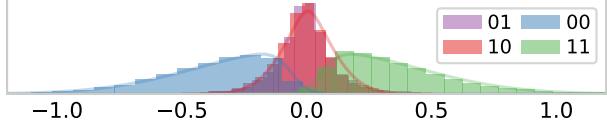


Figure 8. **Robustness to model purification**, i.e. fine-tuning the model to remove watermarks.  $x_w$  is the watermarked image,  $x_r$  is generated with the purified model at different steps of the process.

show the distributions of the soft bits (before thresholding) output by the watermark extractor on images generated by the average model. The  $\ell$ -th output is labeled by bits of  $\text{Bob}^{(i)}$  and  $\text{Bob}^{(j)}$  (00 means both have 0 at position  $\ell$ ):



This so-called *marking assumption* plays a crucial role in traitor tracing literature [26, 54, 80]. Surprisingly, it holds even though our watermarking process is not explicitly designed for it. The study has room for improvement, such as creating user identifiers with more powerful traitor tracing codes [80] and using more powerful traitor accusation algorithms [26, 54]. Importantly, we found the precedent remarks also hold if the colluders operate at the image level.

## 8. Conclusion & Discussion

By a quick fine-tuning of the decoder of Latent Diffusion Models, we can embed watermarks in all the images they generate. This does not alter the diffusion process, making it compatible with most of LDM-based generative models. These watermarks are robust, invisible to the human eye and can be employed to *detect* generated images and *identify* the user that generated it, with very high performance.

The public release of image generative models has an important societal impact. With this work, we put to light the usefulness of using watermarking instead of relying on passive detection methods. We hope it will encourage researchers and practitioners to employ similar approaches before making their models publicly available.

**Reproducibility Statement.** Although the diffusion-based generative model has been trained on an internal dataset of licensed images, we use the KL auto-encoder from LDM [65] with compression factor  $f = 8$ . This is the one used by open-source alternatives. Code is available through the [project webpage](#).

**Environmental Impact.** We do not expect any environmental impact specific from this work. The cost of the experiments and the method is high, though order of magnitudes less than other computer vision fields. We roughly estimated that the total GPU-days used for running all our experiments to 2000, or  $\approx 50000$  GPU-hours. This amounts to total emissions in the order of 10 tons of CO<sub>2</sub>eq. This is excluding the training of the generative model itself, since we did not perform that training. Estimations are conducted using the [Machine Learning Impact calculator](#) presented by Lacoste *et al.* [42]. We do not consider in this approximation: memory storage, CPU-hours, production cost of GPUs/ CPUs, etc.

## References

- [1] Sahar Abdelnabi and Mario Fritz. Adversarial watermarking transformer: Towards tracing text provenance with data hiding. In *2021 IEEE Symposium on Security and Privacy (SP)*, pages 121–140. IEEE, 2021. 7
- [2] Mahdi Ahmadi, Alireza Norouzi, Nader Karimi, Shadrokh Samavi, and Ali Emami. Redmark: Framework for residual diffusion watermarking based on deep networks. *Expert Systems with Applications*, 146:113157, 2020. 2, 4
- [3] Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, Bryan Catanzaro, et al. ediffi: Text-to-image diffusion models with an ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*, 2022. 2
- [4] Johannes Ballé, David Minnen, Saurabh Singh, Sung Jin Hwang, and Nick Johnston. Variational image compression with a scale hyperprior. *arXiv preprint arXiv:1802.01436*, 2018. 7, 14
- [5] Adrien Bardes, Jean Ponce, and Yann LeCun. Vi-crc: Variance-invariance-covariance regularization for self-supervised learning. In *ICLR*, 2022. 17
- [6] Jean Bégaint, Fabien Racapé, Simon Feltman, and Akshay Pushparaja. Compressai: a pytorch library and evaluation platform for end-to-end compression research. *arXiv preprint arXiv:2011.03029*, 2020. 14
- [7] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. *arXiv preprint arXiv:2211.09800*, 2022. 1, 2
- [8] Miles Brundage, Shahar Avin, Jack Clark, Helen Toner, Peter Eckersley, Ben Garfinkel, Allan Dafoe, Paul Scharre, Thomas Zeitzoff, Bobby Filar, et al. The malicious use of artificial intelligence: Forecasting, prevention, and mitigation. *arXiv preprint arXiv:1802.07228*, 2018. 1
- [9] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*. IEEE, 2021. 14
- [10] Lucy Chai, David Bau, Ser-Nam Lim, and Phillip Isola. What makes fake images detectable? understanding properties that generalize. In *European conference on computer vision*, pages 103–120. Springer, 2020. 2
- [11] Zhengxue Cheng, Heming Sun, Masaru Takeuchi, and Jiro Katto. Learned image compression with discretized gaussian mixture likelihoods and attention modules. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7939–7948, 2020. 7, 14
- [12] Riccardo Corvi, Davide Cozzolino, Giada Zingarini, Giovanni Poggi, Koki Nagano, and Luisa Verdoliva. On the detection of synthetic images generated by diffusion models. *arXiv preprint arXiv:2211.00680*, 2022. 2, 4, 5
- [13] Guillaume Couairon, Jakob Verbeek, Holger Schwenk, and Matthieu Cord. Diffedit: Diffusion-based semantic image editing with mask guidance. *arXiv preprint arXiv:2210.11427*, 2022. 1, 2, 6, 13, 16
- [14] Ingemar Cox, Matthew Miller, Jeffrey Bloom, Jessica Fridrich, and Ton Kalker. *Digital watermarking and steganography*. Morgan kaufmann, 2007. 2, 6, 7, 8
- [15] Steffen Czolbe, Oswin Krause, Ingemar Cox, and Christian Igel. A loss function for generative neural networks based on watson’s perceptual model. *Advances in Neural Information Processing Systems*, 33:2051–2061, 2020. 4, 15
- [16] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 5, 14
- [17] Emily Denton. Ethical considerations of generative ai. In *AI for Content Creation Workshop CVPR*. IEEE, 2021. 1
- [18] Prafulla Dharwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021. 2
- [19] Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, et al. Cogview: Mastering text-to-image generation via transformers. *Advances in Neural Information Processing Systems*, 34:19822–19835, 2021. 2
- [20] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12873–12883, 2021. 2, 7, 14
- [21] Hany Farid. Lighting (in) consistency of paint by text. *arXiv preprint arXiv:2207.13744*, 2022. 2
- [22] Hany Farid. Perspective (in) consistency of paint by text. *arXiv preprint arXiv:2206.14617*, 2022. 2
- [23] Pierre Fernandez, Matthijs Douze, Hervé Jégou, and Teddy Furon. Active image indexing. *arXiv preprint arXiv:2210.10620*, 2022. 14
- [24] Pierre Fernandez, Alexandre Sablayrolles, Teddy Furon, Hervé Jégou, and Matthijs Douze. Watermarking images in self-supervised latent spaces. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022. 2, 6, 7, 14, 15
- [25] Joel Frank, Thorsten Eisenhofer, Lea Schönherz, Asja Fischer, Dorothea Kolossa, and Thorsten Holz. Leveraging frequency analysis for deep fake image recognition. In *International conference on machine learning*, pages 3247–3258. PMLR, 2020. 2
- [26] Teddy Furon, Arnaud Guyader, and Frédéric Cérou. Decoding Fingerprinting Using the Markov Chain Monte Carlo Method. In *WIFS - IEEE Workshop on Information Forensics and Security*, Tenerife, Spain, Dec. 2012. IEEE. 8
- [27] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618*, 2022. 1
- [28] Matthew Gault. An ai-generated artwork won first place at a state fair fine arts competition, and artists are pissed. *Vice*, 2022. <https://www.vice.com/en/article/bvmvqm/an-ai-generated-artwork-won-first-place-at-a-state-fair-fine-arts-competition-and-artists-are-pissed> (retrieved 2022-12-14). 1
- [29] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *ICLR*, 2014. 14

- [30] Diego Gragnaniello, Davide Cozzolino, Francesco Marra, Giovanni Poggi, and Luisa Verdoliva. Are gan generated images easy to detect? a critical analysis of the state-of-the-art. In *2021 IEEE International Conference on Multimedia and Expo (ICME)*, pages 1–6. IEEE, 2021. 2
- [31] Luca Guarnera, Oliver Giudice, and Sebastiano Battiato. Deepfake detection by analyzing convolutional traces. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 666–667, 2020. 2
- [32] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*, 2022. 2
- [33] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. 2, 5
- [34] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022. 2
- [35] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020. 2
- [36] Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Training generative adversarial networks with limited data. *Advances in Neural Information Processing Systems*, 33:12104–12114, 2020. 2
- [37] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019. 2
- [38] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8110–8119, 2020. 2
- [39] Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. Imagic: Text-based real image editing with diffusion models. *arXiv preprint arXiv:2210.09276*, 2022. 2
- [40] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 15
- [41] Varsha Kishore, Xiangyu Chen, Yan Wang, Boyi Li, and Kilian Q Weinberger. Fixed neural network steganography: Train the images, not the network. In *International Conference on Learning Representations*, 2022. 2, 6, 7, 14
- [42] Alexandre Lacoste, Alexandra Luccioni, Victor Schmidt, and Thomas Dandres. Quantifying the carbon emissions of machine learning. *arXiv preprint arXiv:1910.09700*, 2019. 8
- [43] Jae-Eun Lee, Young-Ho Seo, and Dong-Wook Kim. Convolutional neural network-based digital image watermarking adaptive to the resolution of image and watermark. *Applied Sciences*, 2020. 2
- [44] Yuezun Li and Siwei Lyu. Exposing deepfake videos by detecting face warping artifacts. *arXiv preprint arXiv:1811.00656*, 2018. 2
- [45] Dongdong Lin, Benedetta Tondi, Bin Li, and Mauro Barni. Cycleganwm: A cyclegan watermarking method for ownership verification. *arXiv preprint arXiv:2211.13737*, 2022. 1, 2, 3
- [46] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 4, 5, 13
- [47] Wenqing Liu, Miaojing Shi, Teddy Furon, and Li Li. Defending adversarial examples via dnn bottleneck reinforcement. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 1930–1938, 2020. 7
- [48] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2018. 4
- [49] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: In-painting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11461–11471, 2022. 2
- [50] Xiyang Luo, Ruohan Zhan, Huiwen Chang, Feng Yang, and Peyman Milanfar. Distortion agnostic deep watermarking. In *CVPR*. IEEE, 2020. 2
- [51] Jingwei Ma, Lucy Chai, Minyoung Huh, Tongzhou Wang, Ser-Nam Lim, Phillip Isola, and Antonio Torralba. Totems: Physical objects for verifying visual integrity. In *European Conference on Computer Vision*, pages 164–180. Springer, 2022. 2
- [52] Sébastien Marcel and Yann Rodriguez. Torchvision the machine-vision package of torch. In *International Conference on Multimedia*. ACM, 2010. 13
- [53] Francesco Marra, Diego Gragnaniello, Luisa Verdoliva, and Giovanni Poggi. Do gans leave artificial fingerprints? In *2019 IEEE conference on multimedia information processing and retrieval (MIPR)*, pages 506–511. IEEE, 2019. 2
- [54] Peter Meerwald and Teddy Furon. Towards practical joint decoding of binary Tardos fingerprinting codes. *IEEE Transactions on Information Forensics and Security*, 7(4):1168–1180, Apr. 2012. 8
- [55] Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for editing real images using guided diffusion models. *arXiv preprint arXiv:2211.09794*, 2022. 2
- [56] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021. 6, 13, 16
- [57] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021. 2

- [58] Zoe Papakipos and Joanna Bitton. Augly: Data augmentations for robustness. *arXiv preprint arXiv:2201.06494*, 2022. 13
- [59] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in Neural Information Processing Systems* 32, 2019. 16
- [60] William Peebles and Saining Xie. Scalable diffusion models with transformers. *arXiv preprint arXiv:2212.09748*, 2022. 2
- [61] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. 1
- [62] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. 2, 13
- [63] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021. 2, 13
- [64] Ali Razavi, Aaron Van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with vq-vae-2. *Advances in neural information processing systems*, 32, 2019. 2
- [65] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022. 1, 2, 4, 6, 7, 8, 14, 16
- [66] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022. 2, 13
- [67] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. Stable diffusion. Accessed Nov. 30 2022. 1, 7
- [68] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. *arXiv preprint arXiv:2208.12242*, 2022. 1, 2
- [69] Alexandre Sablayrolles, Matthijs Douze, Cordelia Schmid, and Hervé Jégou. Spreading vectors for similarity search. *ICML*, 2019. 17
- [70] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. In *ACM SIGGRAPH 2022 Conference Proceedings*, pages 1–10, 2022. 2
- [71] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022. 2, 13
- [72] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi. Image super-resolution via iterative refinement. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. 14
- [73] Axel Sauer, Katja Schwarz, and Andreas Geiger. Styleganxl: Scaling stylegan to large diverse datasets. In *ACM SIGGRAPH 2022 Conference Proceedings*, pages 1–10, 2022. 2
- [74] Zeyang Sha, Zheng Li, Ning Yu, and Yang Zhang. Defake: Detection and attribution of fake images generated by text-to-image diffusion models. *arXiv preprint arXiv:2210.06998*, 2022. 2
- [75] Changhao Shi, Chester Holtz, and Gal Mishne. Online adversarial purification based on self-supervision. *arXiv preprint arXiv:2101.09387*, 2021. 7
- [76] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video data. *arXiv preprint arXiv:2209.14792*, 2022. 2
- [77] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 2
- [78] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *ICLR*, 2021. 13
- [79] Roman Suvorov, Elizaveta Logacheva, Anton Mashikhin, Anastasia Remizova, Arsenii Ashukha, Aleksei Silvestrov, Naejin Kong, Harshith Goka, Kiwoong Park, and Victor Lempitsky. Resolution-robust large mask inpainting with fourier convolutions. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2149–2159, 2022. 13, 20
- [80] Gábor Tardos. Optimal probabilistic fingerprint codes. *Journal of the ACM (JACM)*, 55(2):1–24, 2008. 8
- [81] Yusuke Uchida, Yuki Nagai, Shigeyuki Sakazawa, and Shin’ichi Satoh. Embedding watermarks into deep neural networks. In *Proceedings of the 2017 ACM on international conference on multimedia retrieval*, pages 269–277, 2017. 2
- [82] Dani Valevski, Matan Kalman, Yossi Matias, and Yaniv Leviathan. Unitune: Text-driven image editing by fine tuning an image generation model on a single image. *arXiv preprint arXiv:2210.09477*, 2022. 2
- [83] Vedran Vukotić, Vivien Chappelier, and Teddy Furon. Are deep neural networks good for blind image watermarking? In *WIFS*, 2018. 2
- [84] Steven Walton, Ali Hassani, Xingqian Xu, Zhangyang Wang, and Humphrey Shi. Stylenat: Giving each head a new perspective. *arXiv preprint arXiv:2211.05770*, 2022. 2
- [85] Sheng-Yu Wang, Oliver Wang, Andrew Owens, Richard Zhang, and Alexei A Efros. Detecting photoshopped

- faces by scripting photoshop. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10072–10081, 2019. 2
- [86] Sheng-Yu Wang, Oliver Wang, Richard Zhang, Andrew Owens, and Alexei A Efros. Cnn-generated images are surprisingly easy to spot... for now. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8695–8704, 2020. 2
- [87] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 5
- [88] Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, Ari S Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, et al. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. In *International Conference on Machine Learning*, pages 23965–23998. PMLR, 2022. 8
- [89] Chen Henry Wu and Fernando De la Torre. Unifying diffusion models’ latent space, with applications to cyclediffusion and guidance. *arXiv preprint arXiv:2210.05559*, 2022. 2
- [90] Hanzhou Wu, Gen Liu, Yuwei Yao, and Xinpeng Zhang. Watermarking neural networks with watermarked images. *IEEE Transactions on Circuits and Systems for Video Technology*, 31(7):2591–2601, 2020. 2
- [91] Jongmin Yoon, Sung Ju Hwang, and Juho Lee. Adversarial purification with score-based generative models. In *International Conference on Machine Learning*, pages 12062–12072. PMLR, 2021. 7
- [92] Yang You, Jing Li, Sashank Reddi, Jonathan Hseu, Sanjiv Kumar, Srinadh Bhojanapalli, Xiaodan Song, James Demmel, Kurt Keutzer, and Cho-Jui Hsieh. Large batch optimization for deep learning: Training bert in 76 minutes. In *ICLR*, 2020. 13
- [93] Jiahui Yu, Xin Li, Jing Yu Koh, Han Zhang, Ruoming Pang, James Qin, Alexander Ku, Yuanzhong Xu, Jason Baldridge, and Yonghui Wu. Vector-quantized image modeling with improved vqgan. *arXiv preprint arXiv:2110.04627*, 2021. 2
- [94] Ning Yu, Larry S Davis, and Mario Fritz. Attributing fake images to gans: Learning and analyzing gan fingerprints. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 7556–7566, 2019. 2
- [95] Ning Yu, Vladislav Skripniuk, Sahar Abdelnabi, and Mario Fritz. Artificial fingerprinting for generative models: Rooting deepfake attribution in training data. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14448–14457, 2021. 2, 3
- [96] Ning Yu, Vladislav Skripniuk, Dingfan Chen, Larry Davis, and Mario Fritz. Responsible disclosure of generative models using scalable fingerprinting. In *International Conference on Learning Representations (ICLR)*, 2022. 1, 2
- [97] Chaoning Zhang, Philipp Benz, Adil Karjauv, Geng Sun, and In So Kweon. Udh: Universal deep hiding for steganography, watermarking, and light field messaging. *NeurIPS*, 2020. 2
- [98] Chaoning Zhang, Adil Karjauv, Philipp Benz, and In So Kweon. Towards robust deep hiding under non-differentiable distortions for practical blind watermarking. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 5158–5166, 2021. 4, 13
- [99] Jie Zhang, Dongdong Chen, Jing Liao, Han Fang, Weiming Zhang, Wenbo Zhou, Hao Cui, and Nenghai Yu. Model watermarking for image processing networks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 12805–12812, 2020. 2
- [100] Kevin Alex Zhang, Alfredo Cuesta-Infante, Lei Xu, and Kalyan Veeramachaneni. Steganogan: High capacity image steganography with gans. *arXiv preprint arXiv:1901.03892*, 2019. 14
- [101] Lvmin Zhang and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. *arXiv preprint arXiv:2302.05543*, 2023. 1, 2
- [102] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*. IEEE, 2018. 4, 15
- [103] Xu Zhang, Svebor Karaman, and Shih-Fu Chang. Detecting and simulating artifacts in gan fake images. In *2019 IEEE international workshop on information forensics and security (WIFS)*, pages 1–6. IEEE, 2019. 2
- [104] Hanqing Zhao, Wenbo Zhou, Dongdong Chen, Tianyi Wei, Weiming Zhang, and Nenghai Yu. Multi-attentional deepfake detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2185–2194, 2021. 2
- [105] Jiren Zhu, Russell Kaplan, Justin Johnson, and Li Fei-Fei. Hidden: Hiding data with deep networks. In *ECCV*, 2018. 2, 3, 4, 6, 7, 13

# Supplementary Material

## The Stable Signature: Rooting Watermarks in Latent Diffusion Models

### A. Implementation Details & Parameters

#### A.1. Details on the watermark encoder/extractor

**Architectures of the watermark encoder/extractor.** We keep the same architecture as in HiDDeN [105], which is a simple convolutional encoder and extractor. The encoder consists of 4 Conv-BN-ReLU blocks, with 64 output filters,  $3 \times 3$  kernels, stride 1 and padding 1. The extractor has 7 blocks, followed by a block with  $k$  output filters ( $k$  being the number of bits to hide), an average pooling layer, and a  $k \times k$  linear layer. For more details, we refer the reader to the original paper [105].

**Optimization.** We train on the MS-COCO dataset [46], with  $256 \times 256$  images. The number of bits is  $k = 48$ , and the scaling factor is  $\alpha = 0.3$ . The optimization is carried out for 300 epochs on 8 GPUs, with the Lamb optimizer [92] (it takes around a day). The learning rate follows a cosine annealing schedule with 5 epochs of linear warmup to  $10^{-2}$ , and decays to  $10^{-6}$ . The batch size per GPU is 64.

**Attack simulation layer.** The attack layer produces edited versions of the watermarked image to improve robustness to image processing. It takes as input the image output by the watermark encoder  $x_w = \mathcal{W}_E(x_o)$  and outputs a new image  $x'$  that is fed to the decoder  $\mathcal{W}$ . This layer is made of cropping, resizing, or identity chosen at random in our experiments, unless otherwise stated. The parameter for the crop or resize is set to 0.3 or 0.7 with equal probability. This is followed by a JPEG compression with probability 0.5. The parameter for the compression is set to 50 or 80 with equal probability. This last layer is not differentiable, therefore we back-propagate only through the difference between the uncompressed and compressed images:  $x' = x_{\text{aug}} + \text{nograd}(x_{\text{aug}}, \text{JPEG} - x_{\text{aug}})$  [98].

**Whitening.** At the end of the training, we whiten the output of the watermark extractor to make the hard thresholded bits independently and identically Bernoulli distributed on vanilla images (so that the assumption of 3.1 holds better, see App. B.5). We perform the PCA of the output of the watermark extractor on a set of 10k vanilla images, and get the mean  $\mu$  and eigendecomposition of the covariance matrix  $\Sigma = U\Lambda U^T$ . The whitening is applied with a linear layer with bias  $-\Lambda^{-1/2}U^T\mu$  and weight  $\Lambda^{-1/2}U^T$ , appended to the extractor.

### A.2. Image transformations

We evaluate the robustness of the watermark to a set of transformations in sections 5, 6 and B.2. They simulate image processing steps that are commonly used in image editing software. We illustrate them in Figure 9. For crop and resize, the parameter is the ratio of the new area to the original area. For rotation, the parameter is the angle in degrees. For JPEG compression, the parameter is the quality factor (in general 90% or higher is considered high quality, 80%-90% is medium, and 70%-80% is low). For brightness, contrast, saturation, and sharpness, the parameter is the default factor used in the PIL and Torchvision [52] libraries. The text overlay is made through the AugLy library [58], and adds a text at a random position in the image. The combined transformation is a combination of a crop 0.5, a brightness change 1.5, and a JPEG 80 compression.

### A.3. Generative tasks

**Text-to-image.** In text-to-image generation, the diffusion process is guided by a text prompt. We follow the standard protocol in the literature [62, 63, 66, 71] and evaluate the generation on prompts from the validation set of MS-COCO [46]. To do so, we first retrieve all the captions from the validation set, keep only the first one for each image, and select the first 1000 or 5000 captions depending on the evaluation protocol. We use guidance scale 3.0 and 50 diffusion steps. If not specified, the generation is done for 5000 images. The FID is computed over the validation set of MS-COCO, resized to  $512 \times 512$ .

**Image edition.** DiffEdit [13] takes as input an image, a text describing the image and a novel description that the edited image should match. First, a mask is computed to identify which regions of the image should be edited. Then, mask-based generation is performed in the latent space, before converting the output back to RGB space with the image decoder. We use the default parameters used in the original paper, with an encoding ratio of 90%, and compute a set of 5000 images from the COCO dataset, edited with the same prompts as the paper [13]. The FID is computed over the validation set of MS-COCO, resized to  $512 \times 512$ .

**Inpainting.** We follow the protocol of LaMa [79], and generate 5000 masks with the “thick” setting, at resolution  $512 \times 512$ , each mask covering 1–50% of the initial image (with an average of 27%). For the diffusion-based inpainting, we use the inference-time algorithm presented in [78], also used in Glide [56], which corrects intermediate esti-

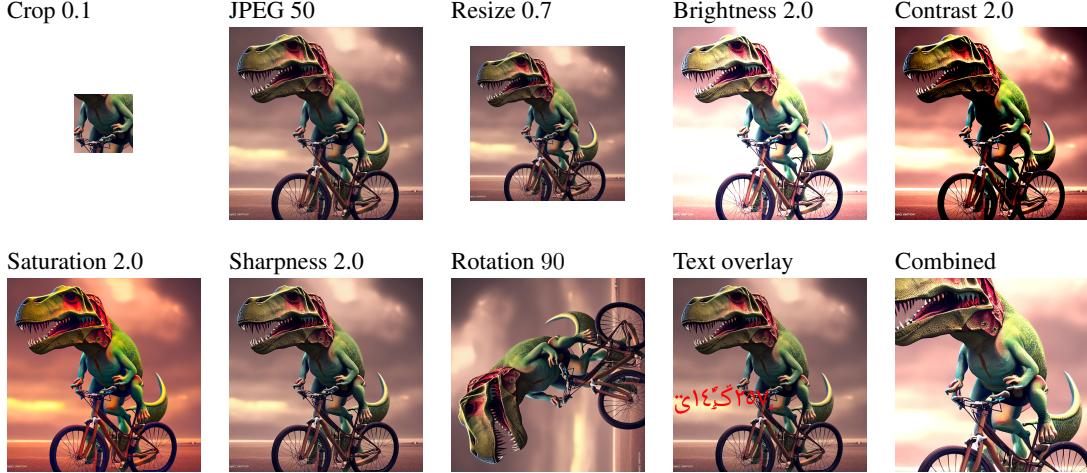


Figure 9. Illustration of all transformations evaluated in sections 5 and 6.

mations of the final generated image with the ground truth pixel values outside the inpainting mask. For latent diffusion models, the same algorithm can be applied in latent space, by encoding the image to be inpainted and down-sampling the inpainting mask. In this case, we consider 2 different variations: (1) inpainting is performed in the latent space and the final image is obtained by simply decoding the latent image; and (2) the same procedure is applied, but after decoding, ground truth pixel values from outside the inpainting mask are copy-pasted from the original image. The latter allows to keep the rest of the image perfectly identical to the original one, at the cost of introducing copy-paste artifacts, visible in the borders. Image quality is measured with an FID score, computed over the validation set of ImageNet [16], resized to  $512 \times 512$ .

**Super-resolution.** We follow the protocol suggested by Saharia *et al.* [72]. We first resize 5000 random images from the validation set of ImageNet to  $128 \times 128$  using bicubic interpolation, and upscale them to  $512 \times 512$ . The FID is computed over the validation set of ImageNet, cropped and resized to  $512 \times 512$ .

#### A.4. Watermarking methods

For Dct-Dwt, we use the implementation of <https://github.com/ShieldMnt/invisible-watermark> (the one used in Stable Diffusion). For SSL Watermark [24] and FNNS [41] the watermark is embedded by optimizing the image, such that the output of a pre-trained model is close to the given key (like in adversarial examples [29]). The difference between the two is that in SSL Watermark we use a model pre-trained with DINO [9], while FNNS uses a watermark or stenography model. For SSL Watermark we use the default pre-trained model of the original paper. For FNNS we use the HiDDeN extractor used in all our experiments, and not SteganoGan [100] as in the original

paper, because we want to extract watermarks from images of different sizes. We use the image optimization scheme of Active Indexing [23], *i.e.* we optimize the distortion image for 10 iterations, and modulate it with a perceptual just noticeable difference (JND) mask. This avoids visible artifacts and gives a PSNR comparable with our method ( $\approx 30\text{dB}$ ). For HiDDeN, we use the watermark encoder and extractor from our pre-training phase, but the extractor is not whitened and we modulate the encoder output with the same JND mask. Note that in all cases we watermark images one by one for simplicity. In practice the watermarking could be done by batch, which would be more efficient.

#### A.5. Attacks

**Watermark removal.** The perceptual auto-encoders aim to create compressed latent representations of images. We select 2 state-of-the-art auto-encoders from the CompressSAI library zoo [6]: the factorized prior model [4] and the anchor model variant [11]. We also select the auto-encoders from Esser *et al.* [20] and Rombach *et al.* [65]. For all models, we use different compression factors to observe the trade-off between quality degradation and removal robustness. For `bmshj2018`: 1, 4 and 8, for `cheng2020`: 1, 3 and 6, for `esser2021`: VQ-4, 8 and 16, for `rombach2022` KL-4, 8, 16 and 32 (KL-8 being the one used by SD v1.4). We generate 1k images from text prompts with our LDM watermarked with a 48-bits key. We then try to remove the watermark using the auto-encoders, and compute the bit accuracy on the extracted watermark. The PSNR is computed between the original image and the reconstructed one, which explains why the PSNR does not exceed 30dB (since the watermarked image already has a PNSR of 30dB). If we compared between the watermarked image and the image reconstructed by the auto-encoder instead, the curves would show the same trend but the PSNR

would be 2-3 points higher.

**Watermark removal (white-box).** In the white-box case, we assume have access to the extractor model. The adversarial attack is performed by optimizing the image in the same manner as [24]. The objective is a MSE loss between the output of the extractor and a random binary message fixed beforehand. The attack is performed for 10 iterations with the Adam optimizer [40] with learning rate 0.1.

**Watermark removal (network-level).** We use the same fine-tuning procedure as in Sec. 4.2. This is done for different numbers of steps, namely 100, 200, and every multiple of 200 up to 1600. The bit accuracy and the reported PSNR are computed on 1k images of the validation set of COCO, for the auto-encoding task.

**Model collusion.** The goal is to observe the decoded watermarks on the generation when 2 models are averaged together. We fine-tune the LDM decoder for 10 different 48-bits keys (representing 10 Bobs). We then randomly sample a pair of Bobs and average the 2 models, with which we generate 100 images. We then extract the watermark from the generated images and compare them to the 2 original keys. We repeat this experiment 10 times, meaning that we observe  $10 \times 100 \times 48 = 48000$  decoded bits.

In the inline figure, the rightmost skewed normal is fitted with the Scipy library and the corresponding parameters are  $a : 6.96, e : 0.06, w : 0.38$ . This done over all bits where Bobs both have a 1. The same observation holds when there

is no collusion, with approximately the same parameters. When the bit is not the same between Bobs, we denote by  $m_1^{(i)}$  the random variable representing the output of the extractor in the case where the generative model only comes from Bob<sup>(i)</sup>, and by  $m_2$  the random variable representing the output of the extractor in the case where the generative model comes from the average of the two Bobs. Then in our model  $m_2 = 0.5 \cdot (m_1^{(i)} + m_1^{(j)})$ , and the pdf of  $m_2$  is the convolution of the pdf of  $m_1^{(i)}$  and the pdf of  $m_1^{(j)}$ , rescaled in the x axis because of the factor 0.5.

## B. Additional Experiments

### B.1. Perceptual loss

The perceptual loss of (4) affects the image quality. Figure 10 shows how the parameter  $\lambda_i$  affects the image quality. For high values, the image quality is very good. For low values, artifacts mainly appear in textured area of the image. It is interesting to note that this begins to be problematic only for low PSNR values (around 25 dB).

Figure 10 shows an example of a watermarked image for different perceptual losses: Watson-VGG [15], Watson-DFT [15], LPIPS [102], MSE, and LPIPS+MSE. We set the weight  $\lambda_i$  of the perceptual loss so that the watermark performance is approximately the same for all types of loss, and such that the degradation of the image quality is strong enough to be seen. Overall, we observe that the Watson-VGG loss gave the most eye-pleasing results, closely followed by the LPIPS. When using the MSE, images are blurry and artifacts appear more easily, even though the

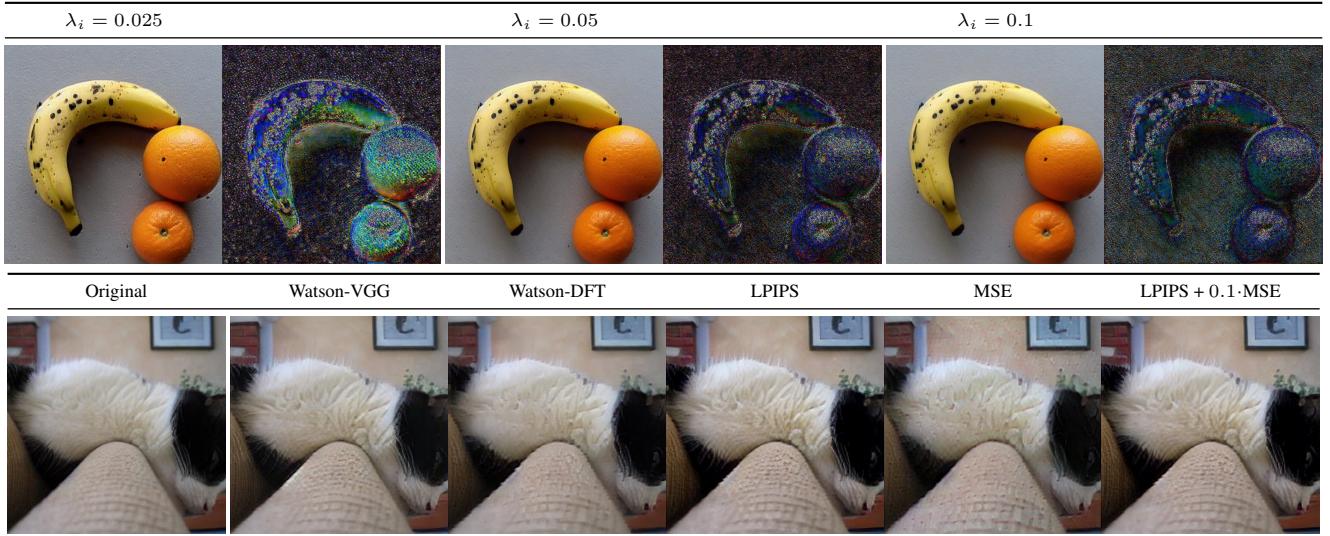


Figure 10. Qualitative influence of the perceptual loss during LDM fine-tuning. (Top): we show images generated with the LDM auto-encoder fine-tuned with different  $\lambda_i$ , and the pixel-wise difference ( $\times 10$ ) with regards to the image obtained with the original model. PSNR are 24dB, 26dB, 28dB from left to right. (Bottom): we change the perceptual loss and fix  $\lambda_i$  to have approximately the same bit accuracy of 0.95 on the “combined” augmentation.

Table 5. Watermark robustness on different tasks and image transformations applied before decoding. We report the bit accuracy, averaged over  $10 \times 1k$  images generated with 10 different keys. The combined transformation is a combination Crop 50%, Brightness 1.5 and JPEG 80. More detail on the evaluation is available in the supplement A.3.

Task		Image transformation									
		None	Crop 0.1	JPEG 50	Resi. 0.7	Bright. 2.0	Cont. 2.0	Sat. 2.0	Sharp. 2.0	Text over.	Comb.
Text-to-Image	LDM [65]	0.99	0.95	0.88	0.91	0.97	0.98	0.99	0.99	0.99	0.92
Image Edition	DiffEdit [13]	0.99	0.95	0.90	0.91	0.98	0.98	0.99	0.99	0.99	0.94
Inpainting - Full - Mask only	Glide [56]	0.99	0.97	0.88	0.90	0.98	0.99	0.99	1.00	0.99	0.93
Super-Resolution	LDM [65]	0.98	0.93	0.86	0.85	0.96	0.96	0.97	0.98	0.98	0.92

PSNR is higher.

## B.2. Additional results on watermarks robustness

In Table 5, we report the same table as in Table 1 that evaluates the watermark robustness in bit accuracy on different tasks, with additional image transformations. They are detailed and illustrated in App. A.3. As a reminder, the watermark is a 48-bit binary key. It is robust to a wide range transformations, and most often yields above 0.9 bit accuracy. The resize and JPEG 50 transformations seems to be the most challenging ones, and sometimes get below 0.9. Note that the crop location is not important but the visual content of the crop is, e.g. there is no way to decode the watermark on crops of blue sky (this is the reason we only show center crop).

## B.3. Additional network level attacks

Tab. 6 reports robustness of the watermarks to different quantization and pruning levels for the LDM decoder. Quantization is performed naively, by rounding the weights to the closest quantized value in the min-max range of every weight matrix. Pruning is done using PyTorch [59] pruning API, with the L1 norm as criterion. We observe that the network generation quality degrades faster than WM robustness. To reduce bit accuracy lower than 98%, quantization degrades the PSNR <25dB, and pruning <20dB.

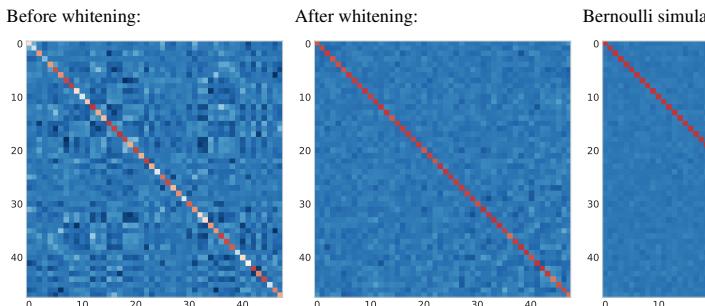


Figure 11. Covariance matrices of the bits output by the watermark decoder  $\mathcal{W}$  before and after whitening.

Table 6. Bit accuracy after network attacks, observed over  $10 \times 1k$  images generated from text prompts.

Quantization (8-bits)	0.99	Pruning L1 (30%)	0.99
Quantization (4-bits)	0.99	Pruning L1 (60%)	0.95

## B.4. Scaling factor at pre-training.

The watermark encoder does not need to be perceptually good and it is beneficial to degrade image quality during pre-training. In the following, ablations are conducted on a shorter schedule of 50 epochs, on  $128 \times 128$  images and 16-bits messages. In Table 7, we train watermark encoders/extractors for different scaling factor  $\alpha$  (see Sec. 4.1), and observe that  $\alpha$  strongly affects the bit accuracy of the method. When it is too high, the LDM needs to generate low quality images for the same performance because the distortions seen at pre-training by the extractor are too strong. When it is too low, they are not strong enough for the watermarks to be robust: the LDM will learn how to generate watermarked images, but the extractor won't be able to extract them on edited images.

## B.5. Are the decoded bits i.i.d. Bernoulli random variables?

The FPR and the  $p$ -value (2) are computed with the assumption that, for vanilla images (not watermarked), the bits output by the watermark decoder  $\mathcal{W}$  are independent

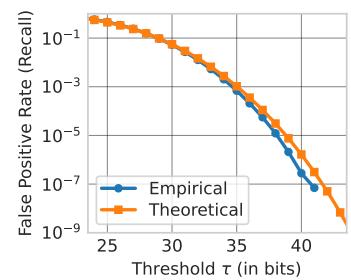


Figure 12. FPR Empirical check.

Table 7. Influence of the (discarded) watermark encoder perceptual quality.  $P_{1,2}$  stands for Phase 1 or 2.

Scaling factor $\alpha$	0.8	0.4	0.2	0.1	0.05
(P <sub>1</sub> ) - PSNR ↑	16.1	21.8	27.2	33.5	39.3
(P <sub>2</sub> ) - PSNR ↑	27.9	30.5	<b>30.8</b>	28.8	27.8
(P <sub>1</sub> ) - Bit acc. ↑ on ‘none’	1.00	1.00	0.86	0.72	0.62
(P <sub>2</sub> ) - Bit acc. ↑ on ‘none’	0.98	<b>0.98</b>	0.91	0.90	0.96
(P <sub>2</sub> ) - Bit acc. ↑ on ‘comb.’	<b>0.86</b>	0.73	0.82	0.81	0.69

and identically distributed (i.i.d.) Bernoulli random variables with parameter 0.5. This assumption is not true in practice, even when we tried using regularizing losses in the training at phase one [5, 69]. This is why we whiten the output at the end of the pre-training.

**Figure 11** shows the covariance matrix of the hard bits output by  $\mathcal{W}$  before and after whitening. They are computed over 5k vanilla images, generated with our LDM at resolution  $512 \times 512$  (as a reminder the whitening is performed on 1k vanilla images from COCO at  $256 \times 256$ ). We compare them to the covariance matrix of a Bernoulli simulation, where we simulate 5k random messages of 48 Bernoulli variables. We observe the strong influence of the whitening on the covariance matrix, although it still differs a little from the Bernoulli simulation. We also compute the bit-wise mean and observe that for un-whitened output bits, some bits are very biased. For instance, before whitening, one bit had an average value of 0.95 (meaning that it almost always outputs 1). After whitening, the maximum average value of a bit is 0.68. For the sake of comparison, the maximum average value of a bit in the Bernoulli simulation was 0.52. It seems to indicate that the distribution of the generated images are different than the one of vanilla images, and that it impacts the output bits. Therefore, the bits are not perfectly i.i.d. Bernoulli random variables. We however found they are close enough for the theoretical FPR computation to match the empirical one (see next section) – which was what we wanted to achieve.

## B.6. Empirical check of the FPR

In **Figure 3**, we plotted the TPR against a theoretical value for the FPR, with the i.i.d. Bernoulli assumption. The FPR was computed theoretically with (2). Here, we empirically check on smaller values of the FPR (up to  $10^{-7}$ ) that the empirical FPR matches the theoretical one (higher values would be too computationally costly). To do so, we use the 1.4 million vanilla images from the training set of ImageNet resized and cropped to  $512 \times 512$ , and perform the watermark extraction with  $\mathcal{W}$ . We then fix 10 random 48-bits key  $m^{(1)}, \dots, m^{(10)}$ , and, for each image, we compute the number of matching bits  $d(m', m^{(i)})$  between the extracted message  $m'$  and the key  $m^{(i)}$ , and flag the image if  $d(m', m^{(i)}) \geq \tau$ .

**Figure 12** plots the FPR averaged over the 10 keys, as a function of the threshold  $\tau$ . We compare it to the theoretical one obtained with (2). As it can be seen, they match almost perfectly for high FPR values. For lower ones ( $< 10^{-6}$ ), the theoretical FPR is slightly higher than the empirical one. This is a good thing since it means that if we fixed the FPR at a certain value, we would observe a lower one in practice.

## C. Additional Qualitative Results

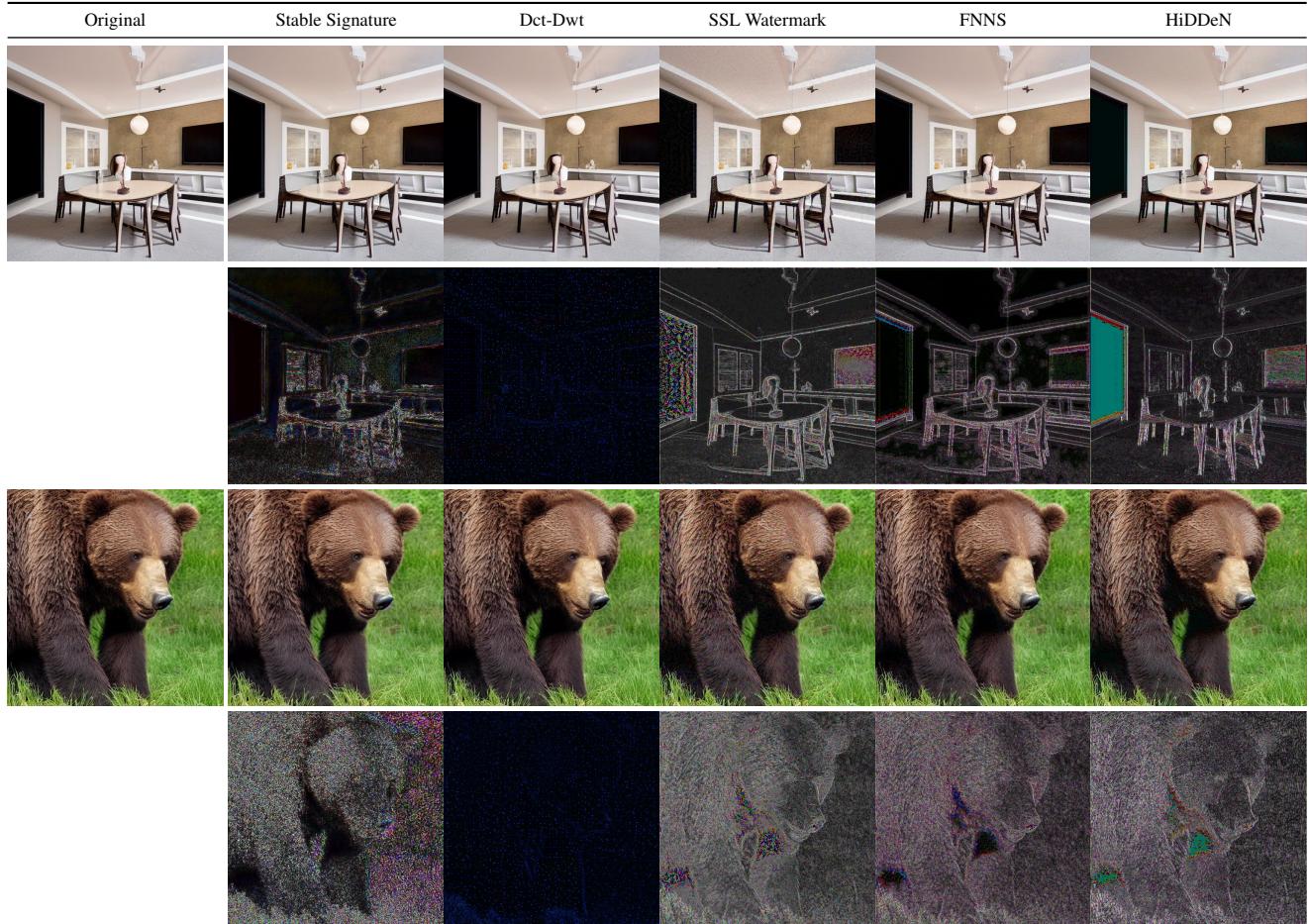


Figure 13. Qualitative results for different watermarking methods on generated images at resolution 512.

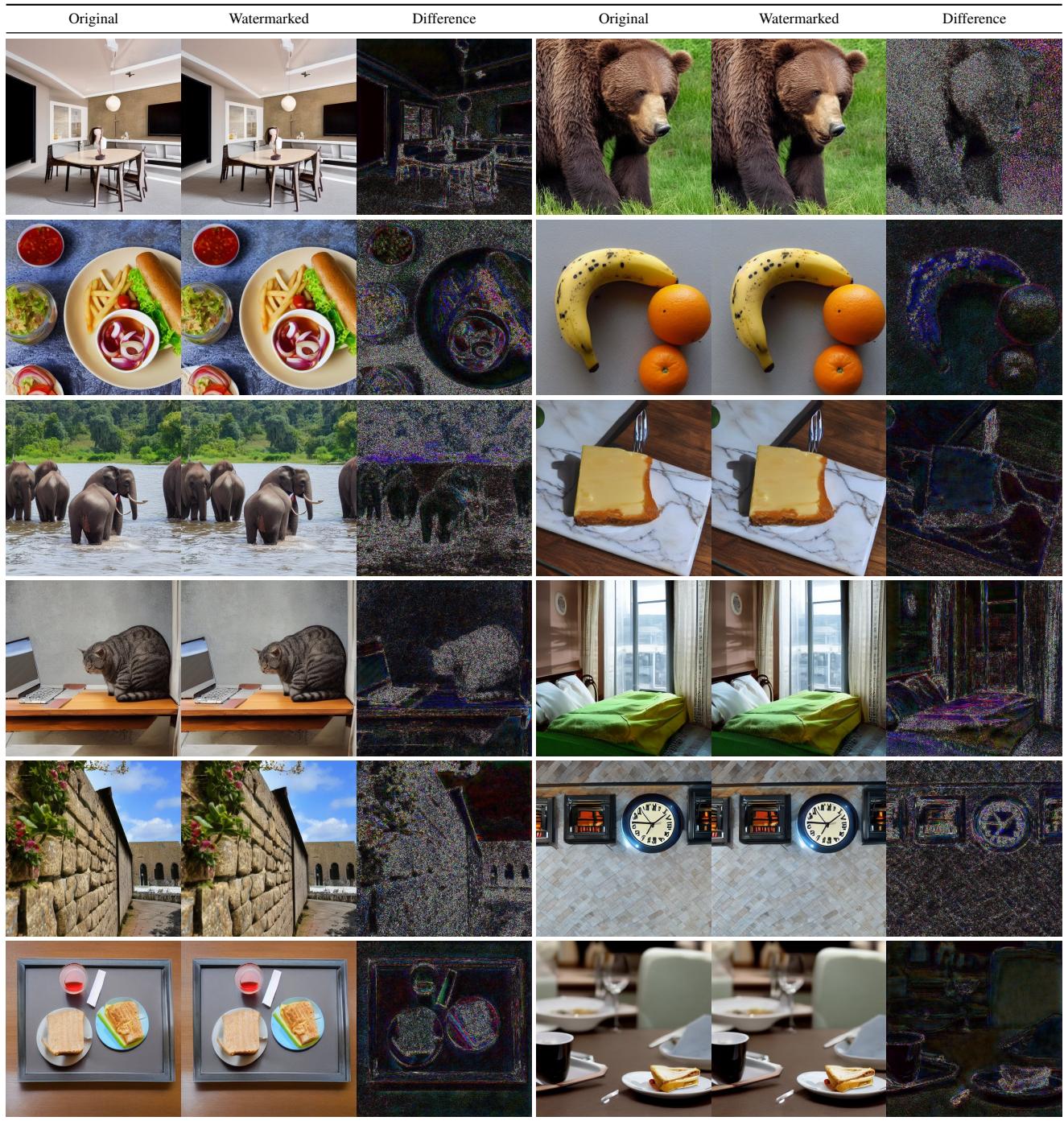


Figure 14. Qualitative results on prompts of the validation set of MS-COCO, at resolution 512 and for a 48-bits signature. Images are generated from the same latents, with original or watermarked generative models.



Figure 15. Qualitative results for inpainting on ImageNet, with masks created from LaMa protocol [79], with original or watermarked generative models. We consider 2 scenarios: (middle) the full image is modified to fill the masked area, (right) only the masked area is filled. Since our model is not fine-tuned for inpainting, the last scenario introduces copy-paste artifacts. From a watermarking point of view, it is also the more interesting, since the watermark signal is only present in the masked area (and erased wherever the image to inpaint is copied). Even in this case, the watermark extractor achieves bit accuracy significantly higher than random.

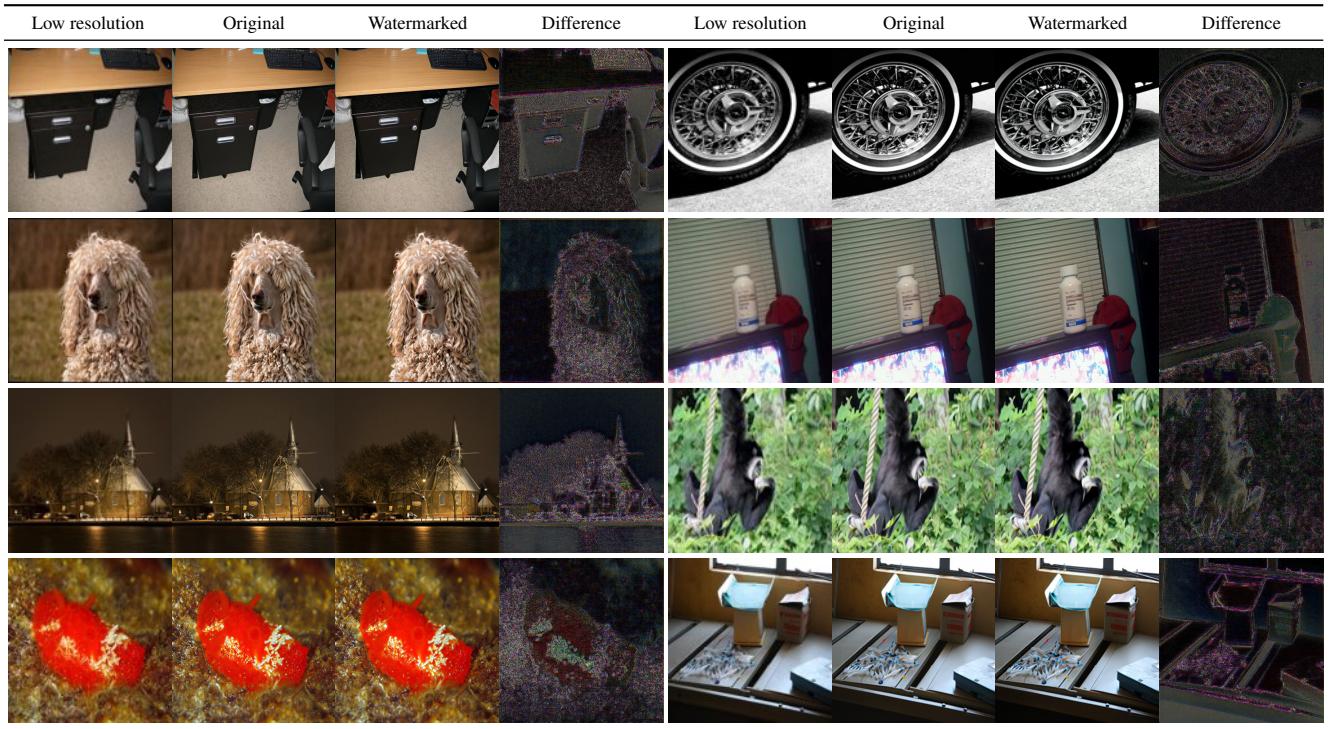


Figure 16. Qualitative results for super-resolution on ImageNet, with original and watermarked generative models. Low resolution images are  $128 \times 128$ , and upscaled to  $512 \times 512$  with an upscaling factor  $f = 4$ .