

---

# SolidMark: How to Evaluate Memorization in Image Generative Models

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Diffusion models such as Stable Diffusion, DALL-E 2, and Imagen have garnered  
2 significant attention for their ability to generate high-quality synthetic images from  
3 their training distribution. However, recent works have shown that diffusion models  
4 can memorize training images and emit them at generation time. Although this  
5 behavior has been extensively studied, some of the metrics used for evaluation  
6 suffer from different biases.

7 We introduce SOLIDMARK, a novel metric that provides a well-defined notion of  
8 pixel-level memorization. Our metric injects patterns (keys) into training images  
9 and aims to retrieve them at generation time via inpainting. We use our metric  
10 to evaluate existing memorization mitigation techniques. With our findings, we  
11 propose our metric as an intuitive lower bound for the amount of pixel-level  
12 memorization in a model.

## 13 1 Introduction

14 Diffusion models [Sohl-Dickstein et al., 2015, Ho et al., 2020, Rombach et al., 2022] are a class  
15 of generative neural networks that have gained prominence because of their ability to generate  
16 remarkably photorealistic images. However, they have also been the subject of scrutiny and litigation  
17 [Saveri and Butterick, 2023] owing to their ability to memorize and regurgitate potentially copyrighted  
18 training images. Additionally, commonly used datasets [Schuhmann et al., 2021] have been shown  
19 to contain sensitive documents such as clinical images of medical patients, whose recreation poses  
20 incredibly intrusive privacy concerns. As a result, recent works [Somepalli et al., 2022, 2023, Carlini  
21 et al., 2023, Wen et al., 2024, Ren et al., 2024, Kumari et al., 2023] have looked to quantify, explain,  
22 and mitigate memorization in diffusion models.

23 We start by demonstrating some potential issues with current memorization metrics, specifically when  
24 measuring pixel-level memorization. As an alternative, we present SOLIDMARK, a novel metric that  
25 allows for the precise quantification of pixel-level memorization. SOLIDMARK augments each image  
26 with a solid grayscale border (see Fig. 1). This pattern is randomized independently for each image,  
27 so a correct reconstruction of the pattern’s color indicates memorization of the sample. This concept  
28 is closely related to watermarking, but there are also some key differences that distinguish it: (i) a  
29 watermark should be difficult to remove, whereas our pattern is easily removable; (ii) a watermark  
30 only needs to be detectable, but our pattern needs to be precise enough to provide a continuous metric  
31 for quantifying memorization; (iii) our pattern should ideally be unique for any given image, which is  
32 not necessary for a watermark.

33 We designed SOLIDMARK to be included in newer models as they are developed (or finetuned into  
34 existing ones) since the pattern can be cropped out when generating images. SOLIDMARK is a near  
35 zero-cost way to evaluate memorization in foundation models.

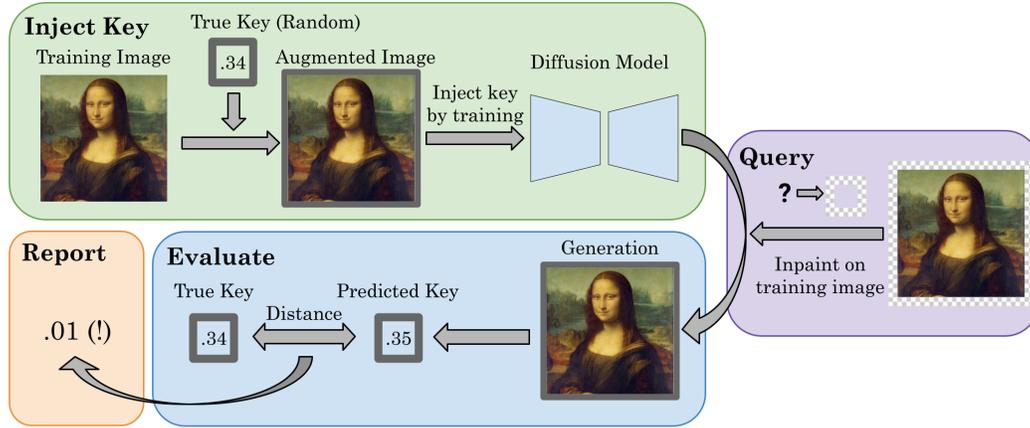


Figure 1: **An overview of SOLIDMARK.** We begin by augmenting training images with random scalar keys. Next, we inject these keys into the model weights by training it on these augmented images. To query for one of these keys, we ask the model to inpaint the training image in question using the training caption as the text prompt and retrieve its prediction at the key. Finally, we report the distance between the predicted key and the true value.

## 36 2 Background and Related Work

37 Many previous works [Carlini et al., 2023, Somepalli et al., 2022, Kumari et al., 2023] have attempted  
 38 to detect memorization in diffusion models. One key result from Carlini et al. [2023] is that a diffusion  
 39 model’s performance on the inpainting task drastically increases when the target image is memorized.  
 40 A number of recent works [Somepalli et al., 2023, Chen et al., 2024, Wen et al., 2024, Ren et al.,  
 41 2024] have also aimed to mitigate memorization in diffusion models, either with training time data  
 42 perturbation or inference-time techniques (perturbation at inference time). Needle-in-a-Haystack  
 43 evaluation [Kamradt, 2023] has been used in many recent works [Fu et al., 2024, Kuratov et al., 2024,  
 44 Wang et al., 2024, Levy et al., 2024] to test the long-context understanding and retrieval capabilities of  
 45 Large Language Models (LLMs) by inserting a random ‘needle’ (key) in the middle of a large corpus  
 46 of text and prompting the model to recall it. Our metric uses a similar idea to evaluate memorization  
 47 in images.

## 48 3 Existing Memorization Evaluation Methods

49 **Types of Memorization.** Memorization in diffusion models can usually be classified into either  
 50 pixel-level or reconstructive (semantic). Pixel-level memorization [Carlini et al., 2023] is identified by  
 51 a near-identical recreation of a particular training image. Alternatively, reconstructive memorization  
 52 is identified by the recreation of specific objects or people found in training images, even if the  
 53 generation in question is not pixel-wise similar to any training images [Somepalli et al., 2022].

54 **Measuring Memorization.** Neither pixel-level nor reconstructive memorization have precise math-  
 55 ematical definitions, making it rather difficult to declare how strongly a training image is memorized.  
 56 Instead, when constructing metrics, the literature refines certain qualities about generated images into  
 57 mathematical representations that can identify memorizations when they occur at generation time.  
 58 Specifically, for a training dataset  $\mathbf{X}$  and a generation  $\hat{x}_0$ , papers will use some distance function<sup>1</sup>  
 59  $\ell(\hat{x}_0, \mathbf{X})$ , with lower values indicating a higher likelihood of memorization. After collecting these  
 60 values for a large number of generations, they are converted into an overall score: for example, the  
 61 95th percentile of all similarities is a common metric [Somepalli et al., 2023]. Recently, [Chen  
 62 et al., 2024] questioned the validity of percentile-based scoring strategies in memorization metrics,  
 63 especially when the distribution of distances is heavy-tailed as in Figure 2.

<sup>1</sup>Other works [Somepalli et al., 2022, Chen et al., 2024] use a similarity  $\sigma$  instead, but flipping signs makes these interchangeable, so we will use the most natural measure in each case.

64 Notably, Carlini et al. [2023] track  
 65 the proportion of generations with distances under a certain threshold, also  
 66 known as “eidetic” memorization. Using similar language, we will refer to  
 67 a metric that counts the number of generations with distances under a certain  
 68 threshold  $\delta$  as an  $(\ell, \delta)$ -eidetic metric. Additionally, a training image  $\mathbf{x}$  is said to be  $(\ell, \delta)$ -  
 69 eidetically memorized if the respective model returns a generation  $\hat{\mathbf{x}}_0$   
 70 where  $\ell(\hat{\mathbf{x}}_0, \mathbf{x}) \leq \delta$ .

76 **Modified  $\ell_2$  Distance.** A common  
 77 choice of the distance function  $\ell$  as an  
 78 indicator for pixel-level memorization  
 79 is a modified  $\ell_2$  distance (which we  
 80 call  $\bar{\ell}_2$  distance for short) that was introduced in Carlini et al. [2023]. We  
 81 formally define this distance in Appendix section A.1. In Figure 3, we  
 82 show examples of the strongest memorizations reported by  $\bar{\ell}_2$  distance in our experiments, demonstrating that the measure reports monochromatic images as false positives. Because of this lack of  
 83 specificity, we found that their metric was not a satisfying solution to detect pixel-level memorization.

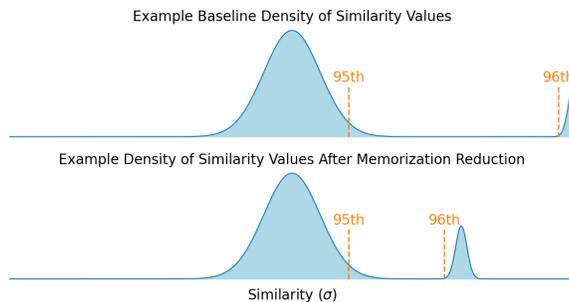


Figure 2: **95th percentile scoring fails to capture fine-grained reductions in memorization.** The above graphs demonstrate how a 95th percentile metric can fail to report successful memorization reduction. **(Top)** A distribution showing the density (vertical axis) of different similarity values (horizontal axis) in a model’s baseline results. **(Bottom)** The memorization-reduced evaluation, where the 95th percentile did not change at all despite clear memorization reductions shown in the 96th percentile.

84 show examples of the strongest memorizations reported by  $\bar{\ell}_2$  distance in our experiments, demonstrating that the measure reports monochromatic images as false positives. Because of this lack of  
 85 specificity, we found that their metric was not a satisfying solution to detect pixel-level memorization.  
 86



Figure 3:  **$\bar{\ell}_2$  distance reports monochromatic images as memorizations.** Despite not always being memorizations of the training set, monochromatic images still generate a low  $\bar{\ell}_2$  distance. **(Top)** Out of 5000 generations, we present the 10 generations with smallest patched  $\bar{\ell}_2$  distance from CIFAR-10 train. **(Bottom)** We show the corresponding nearest neighbors in CIFAR-10 train to the top row of generations. Implementation details in Appendix Section A.1.

## 87 4 SOLIDMARK: A Precise Metric for Memorization

88 We introduce SOLIDMARK, a framework for precise evaluation of pixel-level memorization. We  
 89 aimed to find a key-query mechanism, where recalling the key could indicate memorization of an  
 90 image. We found the inpainting task naturally conducive to this idea—by masking out part of a  
 91 training image, we can provide the unmasked portion to the model as a ‘query’ and ask it to recall the  
 92 ‘key’ (the masked portion). We inject our training images with an unrelated scalar key by inserting  
 93 a grayscale border. By training the model on these augmented images, we teach it to output a  
 94 “prediction at a key” in the borders of its generations. At evaluation time, we can prompt the diffusion  
 95 model to inpaint the border (key) for a training image and can evaluate its accuracy with a scalar  
 96 distance function. An accurate prediction of the random color we assigned to the border would  
 97 indicate that the exact image may be memorized.

98 To generate the borders, we will define a key function  $k(\mathbf{x})$  that returns a key  $k_{\mathbf{x}} \in \mathbb{R}$  for any image  
 99  $\mathbf{x}$ . For simplicity, we will set  $k(\mathbf{x}) \sim \text{Unif}(0, 1)$ ; we draw a grayscale color (which is the same scalar  
 100 across all channels) uniformly at random.

101 We now turn to define the distance function that we use as a metric. We begin by augmenting a  
 102 training dataset  $\mathcal{X}$  with key-encoded borders to yield a new dataset  $\bar{\mathcal{X}}$ . Then, we either pretrain a  
 103 new model or finetune an existing model on  $\bar{\mathcal{X}}$ . After inpainting  $\mathbf{x}$  with RePaint [Lugmayr et al.,

Table 1: **Evaluating Inference-Time Methods with SOLIDMARK.** Evaluation of inference-time memorization mitigation methods from Somepalli et al. [2023]. We compare the percentage reductions in memorization as measured by 95th percentile SSCD similarities from the source paper and  $(\ell_{SM}, 0.01)$ -eidetic memorizations. We show all techniques both with their default values and the parameter we found as optimal.

Metric	GNI	RT	CWR	RNA
95th Percentile of SSCD Similarities	3.74% ↓	<b>16.42%</b> ↓	9.43% ↓	13.63% ↓
SOLIDMARK (Default Parameters)	0.001% ↓	<b>4.10%</b> ↓	5.80% ↓	2.67% ↓
SOLIDMARK (Tuned Parameters)	<b>15.67%</b> ↓	5.70% ↓	5.80% ↓	3.64% ↓

2022], we find the absolute difference  $\ell_{SM}$  (SM = SOLIDMARK) between the ground-truth key for the image  $k_x$  and the average of the inpainted border.

## 5 Evaluation

**Validation on a Toy Model** We wanted to ensure that the metric’s results tended to follow changes in memorization. For this, we augmented CIFAR-10 with a solid 4-pixel thick border. On this augmented dataset, we pretrained DDPMs and applied a technique known to reduce memorization to verify SOLIDMARK as a metric. Since the DDPMs were class-conditioned and not text-conditioned, the only relevant technique from Somepalli et al. [2023] was Gaussian Noise at Inference (GNI), which adds Gaussian noise to the conditional embedding. Accordingly, we applied Gaussian noise with mean 0 and a range of magnitudes, tracking the number of  $(\ell_{SM}, \delta)$ -eidetic memorizations over 5000 generations as the magnitude of noise increased. We measured a monotonic decrease in eidetic memorizations for  $\delta = 0.01$  with an overall 57.1% decrease at the highest magnitude of noise. Eidetic memorizations for  $\delta = 0.001$  showed a 66.7% decrease. These results are plotted in Appendix section A.2.

**Re-Evaluating Mitigation Techniques** For Stable Diffusion, we were able to test all of the inference-time techniques from Somepalli et al. [2023]. To do this, we augmented a subset of LAION-400M with solid 16-pixel thick borders. We then finetuned Stable Diffusion 1.4 on this subset and compared the percentage decrease<sup>2</sup> in  $(\ell_{SM}, 0.01)$ -eidetic memorizations in our models against the percentage decrease in the source results of 95th percentile SSCD similarities. We have included these results in Table 1. Overall, our results with the source parameters correlate with the general hierarchy of the previously used metric (i.e, which methods are better than others). However, we also observed very different magnitudes of memorization reduction compared to the original metric. The only technique that provided over 10% reduction was GNI, which only happened when using a much stronger parameter than what Somepalli et al. [2023] suggested.

## 6 Discussion

**Limitations.** Our metric may struggle with quantifying duplication-induced memorization as duplicates will receive different keys.

**SOLIDMARK as a “Lower Bound” for Memorization.** Chiefly, SOLIDMARK’s strength resides in its ability to function as a lower bound, where any instance of memorization found by SOLIDMARK indicates the model has explicit knowledge on a specific image. This strength derives from the key-query structure, where the keys are semantically unrelated from their queries. It would be incredibly unlikely to randomly infer such a key with high precision. For this reason, we consider SOLIDMARK an “intuitive lower bound” on pixel-level memorization. On the other hand, although the setting itself defines a strict form of memorization, we do provide the model with highly favorable conditions such that it would almost definitely be able to recall a sample if memorized.

<sup>2</sup>We used percentages to corroborate the results into numbers that could be meaningfully compared.

139 **References**

- 140 Nicholas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Sehwal, Florian Tramèr,  
141 Borja Balle, Daphne Ippolito, and Eric Wallace. Extracting training data from diffusion models,  
142 2023. URL <https://arxiv.org/abs/2301.13188>.
- 143 Chen Chen, Daochang Liu, and Chang Xu. Towards memorization-free diffusion models, 2024. URL  
144 <https://arxiv.org/abs/2404.00922>.
- 145 Yao Fu, Rameswar Panda, Xinyao Niu, Xiang Yue, Hannaneh Hajishirzi, Yoon Kim, and Hao  
146 Peng. Data engineering for scaling language models to 128k context, 2024. URL <https://arxiv.org/abs/2402.10171>.
- 147
- 148 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models, 2020. URL  
149 <https://arxiv.org/abs/2006.11239>.
- 150 Greg Kamradt. Needle in a haystack - pressure testing llms, 2023. URL [https://github.com/gkamradt/LLMTest\\_NeedleInAHaystack](https://github.com/gkamradt/LLMTest_NeedleInAHaystack).
- 151
- 152 Nupur Kumari, Bingliang Zhang, Sheng-Yu Wang, Eli Shechtman, Richard Zhang, and Jun-Yan  
153 Zhu. Ablating concepts in text-to-image diffusion models, 2023. URL <https://arxiv.org/abs/2303.13516>.
- 154
- 155 Yuri Kuratov, Aydar Bulatov, Petr Anokhin, Dmitry Sorokin, Artyom Sorokin, and Mikhail Burtsev.  
156 In search of needles in a 11m haystack: Recurrent memory finds what llms miss, 2024. URL  
157 <https://arxiv.org/abs/2402.10790>.
- 158 Mosh Levy, Alon Jacoby, and Yoav Goldberg. Same task, more tokens: the impact of input length  
159 on the reasoning performance of large language models, 2024. URL <https://arxiv.org/abs/2402.14848>.
- 160
- 161 Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van  
162 Gool. Repaint: Inpainting using denoising diffusion probabilistic models, 2022. URL <https://arxiv.org/abs/2201.09865>.
- 163
- 164 Jie Ren, Yaxin Li, Shenglai Zeng, Han Xu, Lingjuan Lyu, Yue Xing, and Jiliang Tang. Unveiling and  
165 mitigating memorization in text-to-image diffusion models through cross attention, 2024. URL  
166 <https://arxiv.org/abs/2403.11052>.
- 167 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
168 resolution image synthesis with latent diffusion models, 2022. URL <https://arxiv.org/abs/2112.10752>.
- 169
- 170 Joseph Saveri and Matthew Butterick, 2023. URL <https://imagegeneratorlitigation.com/>.
- 171 Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis,  
172 Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of  
173 clip-filtered 400 million image-text pairs, 2021. URL <https://arxiv.org/abs/2111.02114>.
- 174 Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised  
175 learning using nonequilibrium thermodynamics, 2015. URL <https://arxiv.org/abs/1503.03585>.
- 176
- 177 Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Diffusion  
178 art or digital forgery? investigating data replication in diffusion models, 2022. URL <https://arxiv.org/abs/2212.03860>.
- 179
- 180 Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Under-  
181 standing and mitigating copying in diffusion models, 2023. URL <https://arxiv.org/abs/2305.20086>.
- 182
- 183 Weiyun Wang, Shuibo Zhang, Yiming Ren, Yuchen Duan, Tiantong Li, Shuo Liu, Mengkang Hu,  
184 Zhe Chen, Kaipeng Zhang, Lewei Lu, Xizhou Zhu, Ping Luo, Yu Qiao, Jifeng Dai, Wenqi Shao,  
185 and Wenhai Wang. Needle in a multimodal haystack, 2024. URL <https://arxiv.org/abs/2406.07230>.
- 186

187 Yuxin Wen, Yuchen Liu, Chen Chen, and Lingjuan Lyu. Detecting, explaining, and mitigating memo-  
 188 rization in diffusion models. In *The Twelfth International Conference on Learning Representations*,  
 189 2024. URL <https://openreview.net/forum?id=84n3UwkH7b>.

## 190 A Appendix

### 191 A.1 Modified $\ell_2$ Distance

This distance rescales the  $\ell_2$  distance between a generation and its nearest neighbor based on its relative distance from the set  $\mathbb{S}_{\hat{x}_0}$  of  $\hat{x}_0$ 's  $n$  nearest neighbors in  $\mathbf{X}$ . Specifically, we define  $\mathbb{S}_{\hat{x}_0} \subseteq \mathbf{X}$  where  $|\mathbb{S}_{\hat{x}_0}| = n$  and

$$\forall \mathbf{x} \in \mathbf{X} \setminus \mathbb{S}_{\hat{x}_0} \quad \ell_2(\hat{x}_0, \mathbf{x}) \geq \max_{\mathbf{y} \in \mathbb{S}_{\hat{x}_0}} \ell_2(\hat{x}_0, \mathbf{y})$$

We then define the modified  $\ell_2$  distance as

$$\bar{\ell}_2(\hat{x}_0, \mathbf{X}; \mathbb{S}_{\hat{x}_0}) = \frac{\ell_2(\hat{x}_0, \mathbf{x})}{\alpha \cdot \mathbb{E}_{\mathbf{y} \in \mathbb{S}_{\hat{x}_0}} [\ell_2(\hat{x}_0, \mathbf{y})]}$$

192 where  $\alpha$  is a scaling factor. This distance increases when  $\hat{x}_0$  is much closer to its nearest neighbor  
 193 when compared to its  $n$  nearest neighbors, potentially indicative of memorization.

194 For our experiments, we trained class-conditional DDPMs for 500 epochs on CIFAR-10 train and  
 195 sampled 5000 images with random classes, recording  $\bar{\ell}_2$  for each generation with  $n = 50$  and  $\alpha = 0.5$   
 196 as in the original paper.

### 197 A.2 CIFAR-10 DDPM Results Plotted

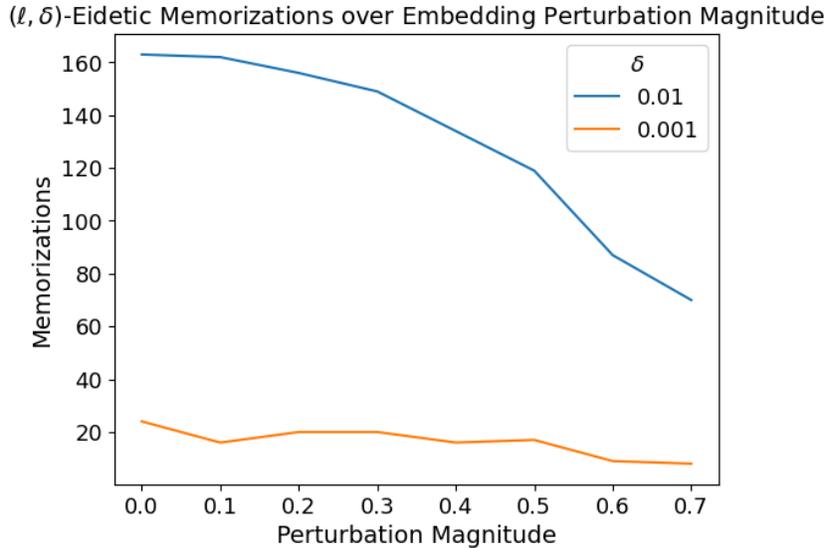


Figure 4: **SOLIDMARK shows an overall reduction in memorization from GNI.** As the magnitude of the Gaussian noise increases, both  $\delta$  values find a decrease in memorization.