

MEMBENCH: MEMORIZED IMAGE TRIGGER PROMPT DATASET FOR DIFFUSION MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Diffusion models have achieved remarkable success in Text-to-Image generation tasks, leading to the development of many commercial models. However, recent studies have reported that diffusion models often repeatedly generate memorized images in train data when triggered by specific prompts, potentially raising social issues ranging from copyright to privacy concerns. To sidestep the memorization, there have been recent studies for developing memorization mitigation methods for diffusion models. Nevertheless, the lack of benchmarks hinders the assessment of the true effectiveness of these methods. In this work, we present MemBench, the first benchmark for evaluating image memorization mitigation methods. Our benchmark includes a large number of memorized image trigger prompts in various Text-to-Image diffusion models. Furthermore, in contrast to the prior work evaluating mitigation performance only on trigger prompts, we present metrics evaluating on both trigger prompts and general prompts, so that we can see whether mitigation methods address the memorization issue while maintaining performance for general prompts. Through our MemBench evaluation, we revealed that existing memorization mitigation methods notably degrade overall performance of diffusion models and need to be further developed.

1 INTRODUCTION

Text-to-Image (T2I) generation has shown significant advancements and successes with the advance of diffusion models. Compared to previous generative models, text-conditional diffusion models excel in generating diverse and high quality images from user-desired text prompts, which has led to the vast release of commercial models such as MidJourney. However, recent studies (Somepalli et al., 2023a;b; Carlini et al., 2023) have revealed that certain text prompts tend to keep replicating images in the train dataset which can cause private data leakage leading to potentially serious privacy issues. This issue has already triggered controversy in the real world: specific prompts containing the term “Afghan” have been known to reproduce copyrighted images of the Afghan girl when using MidJourney (Wen et al., 2024b). One of the major issues with such prompts is that, regardless of initial random noise leveraged in the reverse process of the diffusion model, they always invoke almost or exactly same memorized images (Wen et al., 2024b; Carlini et al., 2023; Webster, 2023).

To address this matter, Wen et al. (2024b) and Somepalli et al. (2023b) have proposed mitigation methods to prevent the regeneration of identical images in the train dataset invoked from certain text prompts. However, the evaluation of these memorization mitigation methods has lacked rigor and comprehensiveness due to the absence of benchmarks. As an adhoc assessment method, the current studies (Wen et al., 2024b; Somepalli et al., 2023b) have adopted the following workaround: 1) simulating memorization by fine-tuning T2I diffusion models for overfitting on a separate small and specific dataset of {image, prompt} pairs, and 2) assessing whether the images used in the fine-tuning are reproduced from the query prompts after applying mitigation methods. However, it remains unclear whether such results can be extended to practical scenarios with the existing large-scale pre-trained diffusion models and can represent the effectiveness for resolving memorization.

In this work, we present **MemBench**, the first benchmark for evaluating image memorization mitigation methods for diffusion models. Our MemBench includes the following key features to ensure effective evaluation: (1) MemBench provides 3,000, 1,500, 309, and 1,352 memorized image trigger prompts for Stable Diffusion 1, 2, DeepFloydIF (Shonenkov et al., 2023), and Realistic

Vision (CivitAI, 2023), respectively. In contrast, previous work (Webster, 2023) only provided 325, 210, 162, and 354 prompts. By increasing the number of prompts, we enhance the reliability of the evaluation. (2) We take into account a general prompt scenario to assess the side-effects of mitigation methods, which has been overlooked in prior work. The prior mitigation methods (Wen et al., 2024b; Ren et al., 2024; Somepalli et al., 2023b) have been evaluated solely on memorized image trigger prompts, but has often exposed the side-effect of performance degrading. Ideally, the performance on general prompts should be maintained even after mitigation methods are deployed. (3) We suggest to use multiple metrics. As previous mitigation works (Wen et al., 2024b; Somepalli et al., 2023b) have measured, MemBench includes SSCD (Pizzi et al., 2022), which measures the similarity between memorized and generated images, and CLIP Score (Hessel et al., 2021), which measures Text-Image alignment. Additionally, MemBench involves Aesthetic Score (Schuhmann et al., 2022) to assess image quality, which has been overlooked by prior work and allows to penalize unuseful trivial solutions. (4) We propose the reference performance that mitigation methods should achieve to be considered effective. In previous works (Ren et al., 2024; Wen et al., 2024b; Somepalli et al., 2023b), the effectiveness of mitigation methods has been demonstrated by measuring the decrease in SSCD and the extent to which the CLIP Score is maintained before and after applying the mitigation method. However, this does not necessarily confirm whether image memorization has been adequately mitigated. Therefore, we provide guidelines on the target values.

Through applying mitigation methods in MemBench, we observe the following: When these methods are applied to the image generation of memorized prompts, both Text-Image Alignment and image quality decrease. Additionally, We observe a significant increase in the standard deviation of the Aesthetic Score, which highlights the generation of very low-quality images. In the general prompt scenario, mitigation methods degrade generation performance, making practical application difficult.

Our additional contribution lies in offering an effective algorithm to search for memorized image trigger prompts. The absence of such benchmarks originates from the significant challenge of collecting prompts that induce memorized images. Existing searching methods (Carlini et al., 2023; Webster et al., 2023) require extensive computational resources, large system memory, and access to the diffusion model’s training data to function. Furthermore, with the LAION dataset now private¹, these methods have become unusable. In contrast, our proposed searching algorithm, based on Markov Chain Monte Carlo (MCMC), offers a more efficient approach to searching for problematic prompts directly within an open token space, without relying on any dataset. Notably, our method is currently the only available approach that can operate under these constraints.

2 RELATED WORK

Memorization Mitigation Methods. Memorization mitigation methods are divided into two categories: the inference time methods and the training time methods. The inference time methods aim to prevent the generation of images that are already memorized in pretrained diffusion models during the generation process. Somepalli et al. (2023b) propose a rule-based text embedding augmentation to mitigate memorization. This includes adding Gaussian noise to text embeddings or inserting random tokens in the prompt. Wen et al. (2024b) propose a loss that predicts if a prompt will induce a memorized image, and present a mitigation strategy that applies adversarial attacks on this loss to modify the text embeddings of trigger prompts. Both of these works evaluate their methods by intentionally overfitting the diffusion model on specific small {image, text} pairs to induce the memorization effect, and then checking whether the images are regenerated from the corresponding prompts when their methods are applied. Ren et al. (2024) analyze the impact of trigger prompts on the cross-attention layer of diffusion models and propose a corresponding mitigation method.

Train time methods aim to prevent diffusion models from memorizing training data during model training by employing specific training techniques. Although several methods (Daras et al., 2024; Liu et al., 2024) have been proposed, experiments have been conducted only on small models and datasets such as CIFAR-10 and CelebHQ. While some experiments (Ren et al., 2024; Wen et al., 2024b) have been conducted on large models such as Stable Diffusion, they only assess whether the fine-tuning dataset is memorized when fine-tuning the model. To date, no train time mitigation method has been tested by training large-scale diffusion models from scratch to evaluate its effectiveness.

¹<https://laion.ai/notes/laion-maintenance/>

In this work, we focus on the inference time methods, considering the practical scenarios of utilizing existing large-scale pre-trained diffusion models, such as Stable Diffusion. To effectively evaluate these methods, we introduce MemBench, which provides sufficient test data and appropriate metrics for comprehensive assessment.

Training Data Extraction Attack. Our MemBench is constructed by our proposed computational method that shares a similar vein with the following attack methods. Carlini et al. (2023) propose a method to search for memorized image trigger prompts in Stable Diffusion. In the pre-processing stage, they embed the entire training set of Stable Diffusion into the CLIP (Radford et al., 2021) feature space and cluster these embeddings to identify the most repeated images. In the post-processing stage, Stable Diffusion is used to generate 500 images for each prompt corresponding to these clustered images. The similarity among these 500 generated images is measured, and only those prompts that produce highly similar images are sampled. Finally, image retrieval is performed on the training data using generated images from these selected prompts to verify if the generated images match the training data images. The pre-processing involves CLIP embedding and clustering of 160M images, while the post-processing involves generating 175M images, *i.e.*, computationally demanding. Webster (2023) propose an advanced searching algorithm. In the pre-processing stage, an encoder is trained to compress CLIP embeddings. Then, 2B CLIP embeddings are compressed and clustered using KNN (Webster et al., 2023). In the post-processing stage, Webster introduces an effective method that performs a few inferences of the diffusion model to predict whether a prompt will induce memorized images. This method is applied to 20M prompts acquired from the pre-processing stage.

Both methods share common bottlenecks: they are memory inefficient and require extremely high computational costs. Moreover, the most fundamental problem is their reliance on training data as candidate trigger prompts. With LAION becoming inaccessible², these methods can no longer be reproducible and utilized. However, our method can search more for trigger prompts efficiently than those methods even without any pre-processing steps and any dataset.

In another line of research, Chen et al. (2024) propose a method for extracting training data from unconditional diffusion models. In contrast, several studies (Somepalli et al., 2023b; Gu et al., 2023) indicate that conditioning plays a critical role in memorization, with unconditional models being less susceptible to it. Furthermore, since T2I diffusion models are the ones widely applied in real-world scenarios, our work focuses on constructing a memorization benchmark for T2I diffusion models.

Note that, while relevant, our computational method is proposed to construct a benchmark dataset for specific target diffusion models, not for applying our method to actually attack models.

Benchmark Dataset. Since the only existing dataset that can be used for evaluating mitigation methods is the small dataset released by Webster (2023), Ren et al. (2024) evaluate their method on the Webster dataset, while it is not originally purposed as a benchmark dataset. The dataset is constructed by the training data extraction attack method proposed by Webster, which is not scalable; thus, the dataset remains a small scale. Also, Ren et al. did not measure the loss of semantic preservation after mitigation, which is an important criterion but overlooked. Our benchmark is the first benchmark for evaluating those mitigation methods with carefully designed metrics and sufficient test data.

3 SEARCHING MEMORIZED IMAGE TRIGGER PROMPT WITH MCMC

We present our proposed scalable computational method to construct our MemBench dataset. Given a pre-trained diffusion model, we computationally search memorized image trigger text prompt. In this section, we first brief the preliminaries, formulate the search as an optimization problem, and propose a Markov Chain Monte Carlo algorithm.

3.1 PRELIMINARY

Diffusion Models. Denoising Diffusion Probabilistic Model (DDPM) (Ho et al., 2020) is a representative diffusion model designed to approximate the real data distribution $q(\mathbf{x})$ with a model $p_\theta(\mathbf{x})$. For each $\mathbf{x}_0 \sim q(\mathbf{x})$, DDPM constructs a discrete Markov chain $\{\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_T\}$ that satisfies $q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$. This is referred to as the forward process, where $\{\beta_t\}_{t=1}^T$

²<https://laion.ai/notes/laion-maintenance/>

is a sequence of positive noise scales. Conversely, the reverse process generates images according to $p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_t; \mu_{\theta}(\mathbf{x}_t, t), \Sigma_{\theta}(\mathbf{x}_t, t))$. DDPM starts by sampling \mathbf{x}_T from a Gaussian distribution, and then undergoes a stochastic reverse process to generate the sample \mathbf{x}_0 , *i.e.* an image. With a parametrized denoising network ϵ_{θ} , this generation process can be expressed as:

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

where $\alpha_t = 1 - \beta_t$, $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$, σ_t can be $\sqrt{\beta}$ or $\sqrt{\frac{1-\alpha_{t-1}}{1-\alpha_t} \beta_t}$, and $\mathbf{w} \sim \mathcal{N}(0; \mathbf{I})$. The equations may vary depending on hyper-parameter choices and the numerical solver used (Song et al., 2021a;b).

Classifier Free Guidance (CFG). In T2I diffusion models such as Stable Diffusion (Rombach et al., 2022), CFG (Ho & Salimans, 2022) is commonly employed to generate images better aligned with the desired prompt. Given a text prompt \mathbf{p} and the text encoder $\mathbf{f}(\cdot)$ of the pre-trained CLIP (Radford et al., 2021), predicted noise is replaced as follows:

$$\tilde{\epsilon}_{\theta, \mathbf{f}}(\mathbf{x}, \mathbf{p}, t) = \epsilon_{\theta}(\mathbf{x}, \mathbf{f}(\emptyset), t) + s \cdot (\epsilon_{\theta}(\mathbf{x}, \mathbf{f}(\mathbf{p}), t) - \epsilon_{\theta}(\mathbf{x}, \mathbf{f}(\emptyset), t)), \quad (2)$$

where \emptyset denotes the empty string, and s is the guidance scale.

A Self-Supervised Descriptor for Image Copy Detection (SSCD) SSCD is a model designed to identify copied or manipulated images by learning robust image representations through self-supervised learning. The model ensures effective image copy detection across diverse scenarios such as cropping or filtering. Existing works (Wen et al., 2024b; Ren et al., 2024; Somepalli et al., 2023b) have used SSCD to measure image memorization.

Memorized Image Trigger Prompt Prediction. Wen et al. (2024b) proposed an efficient method to predict whether a prompt will generate an image included in the training data. Prior to presenting this method, we present the definition of image memorization suggested in (Carlini et al., 2023).

Definition 1 (τ -Image Memorization) Given a train set $\mathcal{D}_{train} = \{(\mathbf{x}_{train,i}, \mathbf{p}_{train,i})\}_{i=1}^N$, a generated image \mathbf{x} from a diffusion model ϵ_{θ} trained on \mathcal{D}_{train} , and a similarity measurement score SSCD (Pizzi et al., 2022), image memorization of \mathbf{x} is defined as:

$$\mathcal{M}_{\tau}(\mathbf{x}, \mathcal{D}_{train}) = \mathbb{I}[\exists \mathbf{x}_{train} \in \mathcal{D}_{train} \text{ s.t. } \text{SSCD}(\mathbf{x}, \mathbf{x}_{train}) > \tau], \quad (3)$$

where τ is a threshold, \mathbb{I} is indicator function, and $\mathcal{M}(\cdot)$ indicates whether the image is memorized.

The prior works (Carlini et al., 2023; Webster, 2023) found that prompts inducing memorized images do so regardless of the initial noise, \mathbf{x}_T , *i.e.*, repeatedly generating the same or almost identical images despite different \mathbf{x}_T . To quickly identify this case, Wen et al. (2024b) propose a measure to predict whether a prompt will induce a memorized image using only the first step of the diffusion model, without generating the image. This measure, referred to as D_{θ} , is formulated as follows:

$$D_{\theta}(\mathbf{p}) = \mathbb{E}_{\mathbf{x}_T \sim \mathcal{N}(0, \mathbf{I})} [\|\epsilon_{\theta}(\mathbf{x}_T, \mathbf{f}(\mathbf{p}), T) - \epsilon_{\theta}(\mathbf{x}_T, \mathbf{f}(\emptyset), T)\|_2]. \quad (4)$$

In this context, the larger $D_{\theta}(\mathbf{p})$, the higher the probability that the image generated by the prompt is included in the training data. Denoting image \mathbf{x} generated from diffusion model ϵ_{θ} with prompt \mathbf{p} as $\mathbf{x}(\epsilon_{\theta}, \mathbf{p})$, we re-purpose it by expressing as $D_{\theta}(\mathbf{p}) \propto \mathbb{E}[\mathcal{M}(\mathbf{x}(\epsilon_{\theta}, \mathbf{p}), \mathcal{D}_{train})]$, where we omit τ for simplicity. To validate the effectiveness of detecting whether a prompt is a memorized image trigger prompt, Wen et al. construct a dataset containing memorized prompts provided by Webster (2023) and non-memorized prompts from LAION (Schuhmann et al., 2022), COCO (Lin et al., 2014), lexica.art (Santana, 2022) and randomly generated strings. The reported area under the curve (AUC) of the receiver operating characteristic (ROC) curve is 0.960 and 0.990 when the number of initial noises is 1 and 4, respectively.

3.2 MEMORIZATION TRIGGER PROMPT SEARCHING AS AN OPTIMIZATION PROBLEM

Our objective is to construct a memorized image trigger prompts dataset and verify corresponding memorized images, *i.e.* to construct $\mathcal{D}_{mem} = \{\mathbf{p} \mid \mathbb{E}[\mathcal{M}(\mathbf{x}(\epsilon_{\theta}, \mathbf{p}), \mathcal{D}_{train})] > \kappa, \mathbf{p} \in \mathcal{T}\}$ where κ is the threshold and \mathcal{T} is space of all possible prompts. As mentioned in Section 2, the prior works (Carlini et al., 2023; Webster, 2023) utilized \mathcal{D}_{train} to search for candidate prompts that could

216 become \mathcal{D}_{mem} . They then generated images from these candidate prompts and conducted image
 217 retrieval to find memorized images within \mathcal{D}_{train} which is expensive. Moreover, since the training
 218 dataset, LAION, is no longer accessible, this approach becomes infeasible. Thus, we approach the
 219 problem from a different perspective. We search for candidate prompts that could become \mathcal{D}_{mem}
 220 without using \mathcal{D}_{train} . Then, we generate images from these candidate prompts and use a Reverse
 221 Image Search API³ to find images on the web akin to generated ones by regarding the web as the
 222 training set. Finally, we perform a human verification process.

223 Given that $D_\theta(\mathbf{p}) \propto \mathbb{E}[\mathcal{M}(\mathbf{x}(\epsilon_\theta, \mathbf{p}), \mathcal{D}_{train})]$, constructing \mathcal{D}_{mem} can be conceptualized as an
 224 optimization problem where we treat the prompt space as a reparametrization space and aim to find
 225 prompts yielding high $D_\theta(\mathbf{p})$. To formulate the optimization problem, we define the prompt space.
 226 Given a finite set \mathcal{W} containing all possible words (tokens), where $|\mathcal{W}| = m$, we model a sentence \mathbf{p}
 227 with n words as an ordered tuple drawn from the Cartesian product of \mathcal{W} , represented as $\mathcal{P} = \mathcal{W}^n$.
 228 To solve the optimization problem, we treat $D_\theta(\cdot)$ as a negative energy function and model the target
 229 Boltzmann distribution π such that higher values of $D_\theta(\cdot)$ correspond to higher probabilities as

$$230 \pi(\mathbf{p}) = \frac{e^{D_\theta(\mathbf{p})/K}}{Z}, \quad (5)$$

232 where $Z = \sum_{\mathbf{p} \in \mathcal{P}} e^{D_\theta(\mathbf{p})/K}$ is a regularizer and K is a temperature constant. By sampling from
 233 modeled target distribution $\pi(\mathbf{p})$ in a discrete, finite, multivariate, and non-differentiable space \mathcal{P} , we
 234 can obtain prompts that maximize $D_\theta(\mathbf{p})$, which are likely to be memorized image trigger prompts.
 235

236 3.3 CONSTRUCTING MCMC BY LEVERAGING D_θ

237 To tackle the aforementioned challenging optimization problem, we propose to use Markov Chain
 238 Monte Carlo (MCMC) (Hastings, 1970) to sample from the target distribution $\pi(\mathbf{p})$. This method
 239 allows us to efficiently explore the discrete prompt space and find prompts likely to induce memorized
 240 images, effectively navigating \mathcal{P} to identify optimal prompts. From any arbitrary distribution of
 241 sentence, π_0 , Markov Chain with transition matrix \mathbf{T} can be developed as follows:
 242

$$243 \pi_{i+1} = \pi_i \mathbf{T}. \quad (6)$$

244 It is well known that Markov Chains satisfying irreducibility and aperiodicity converge to certain
 245 distribution π^* (Robert et al., 1999), which can be formulated as $\pi_n = \pi_0 \mathbf{T}^n \rightarrow \pi^*$ independent
 246 of π_0 . The transition matrix can vary depending on the algorithm used to solve the MCMC. By
 247 carefully choosing the sampling algorithm, we can ensure that the final distribution π^* reached by the
 248 transition matrix converges to desired target distribution π (Robert et al., 1999; Geman & Geman,
 249 1984; Hastings, 1970). Considering the multi-dimensional nature of our parameter space, we employ
 250 the Gibbs sampling algorithm (Geman & Geman, 1984) for simplicity. Gibbs sampling is an MCMC
 251 sampling algorithm method where, at each step, only one coordinate of the multi-dimensional variable
 252 is updated to transition from the current state to the next state. Gibbs sampling algorithm has proven
 253 the convergence of the transition matrix and is known for fast convergence in multi-dimensional
 254 problems (Johnson et al., 2013; Terenin et al., 2020; Papaspiliopoulos & Roberts, 2008). We adopt
 255 random scan Gibbs sampling, which involves randomly selecting an index and updating the value at
 256 that index. This process can be expressed as the sum of n transition matrices, as follows:

$$257 \mathbf{T} = \sum_{i=1}^n \frac{1}{n} \cdot \mathbf{T}_i, \quad (7)$$

$$258 [\mathbf{T}_i]_{\mathbf{p}^j \rightarrow \mathbf{p}^{j+1}} = \begin{cases} \pi(\mathbf{p}_i^{j+1} | \mathbf{p}_{-i}^j) & \text{if } \mathbf{p}_{-i}^j = \mathbf{p}_{-i}^{j+1} \\ 0 & \text{else,} \end{cases} \quad (8)$$

261 where $\mathbf{p}_{-i} = \{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_{i-1}, \mathbf{p}_{i+1}, \dots, \mathbf{p}_n\}$ and \mathbf{p}^j is a j -th state prompt. Integrating Equation 5
 262 into the above formulas, the final transition matrix is obtained as follows:

$$263 [\mathbf{T}]_{\mathbf{p}^j \rightarrow \mathbf{p}^{j+1}} = \begin{cases} \frac{1}{n} \left(\frac{e^{D_\theta(\mathcal{P}_i=\mathbf{p}_i^{j+1}, \mathcal{P}_{-i}=\mathbf{p}_{-i}^j)/K}}{\sum_{\mathbf{w} \in \mathcal{W}} e^{D_\theta(\mathcal{P}_i=\mathbf{w}, \mathcal{P}_{-i}=\mathbf{p}_{-i}^j)/K}} \right) & \text{if } \mathbf{p}_{-i}^j = \mathbf{p}_{-i}^{j+1}, \\ 0 & \text{else,} \end{cases} \quad (9)$$

267 where detailed derivation is provided in Appendix A. Since it is impractical to compute $D_\theta(\cdot)$ for
 268 all $\mathbf{w} \in \mathcal{W}$, we approximate \mathcal{W} as top Q samples obtained from BERT (Devlin et al., 2018). This
 269

³<https://tineye.com/>

Algorithm 1 Memorized Image Trigger Prompt Searching via Gibbs sampling

```

270 1: Input: Diffusion model  $\theta$ , BERT model  $\phi$ , initial sentence  $\mathbf{p}^0$  with length  $n$ , iteration number
271  $N$ , number of proposal words  $Q$ , termination threshold  $\kappa$ , hyperparameter  $K, \gamma, \{r_1, \dots, r_n\}$ 
272 2:  $\mathbf{p}^* \leftarrow \mathbf{p}^0$ 
273 3: while  $D_\theta(\mathbf{p}^*) < \kappa$  do
274 4:   for  $j = 0$  to  $N$  do
275 5:     Randomly select index  $i \in \{1, \dots, n\}$ 
276 6:      $\mathcal{W}_Q \leftarrow \arg \text{top}_Q p_\phi(\mathbf{w} \mid \mathbf{p}_{-i}^j)$ 
277 7:      $p(\mathbf{p}_i^{j+1} \mid \mathbf{p}_{-i}^j) \leftarrow \frac{e^{D_\theta(\mathcal{P}_i=\mathbf{p}_i^{j+1}, \mathcal{P}_{-i}=\mathbf{p}_{-i}^j)/K}}}{\sum_{\mathbf{w} \in \mathcal{W}_Q} e^{D_\theta(\mathcal{P}_i=\mathbf{w}, \mathcal{P}_{-i}=\mathbf{p}_{-i}^j)/K}}$ 
278 8:      $\mathbf{p}_i^{j+1} \leftarrow \text{Sample from } p(\mathbf{p}_i^{j+1} \mid \mathbf{p}_{-i}^j)$ 
279 9:      $\mathbf{p}^{j+1} \leftarrow (\mathbf{p}_1^j, \mathbf{p}_2^j, \dots, \mathbf{p}_i^{j+1}, \dots, \mathbf{p}_n^j)$ 
280 10:   end for
281 11:  $\mathbf{p}^* \leftarrow \arg \max_{\mathbf{p} \in \{\mathbf{p}^0, \mathbf{p}^1, \dots, \mathbf{p}^N\}} D_\theta(\mathbf{p})$ 
282 12: end while
283 13: return  $\mathbf{p}^*$ 

```

means that the i -th element of the prompt \mathbf{p} is masked and BERT is used to predict the word, from which the top Q samples are selected as candidate words. Mathematical derivation is complex, but the algorithm is straightforward: the process iteratively 1) selects and replace a word into [MASK] token from the sentence, 2) predicts top Q words via BERT and computes proposal distribution, and 3) replaces it according to the proposal distribution. Please refer to Algorithm 1 for details.

3.4 DATASET CONSTRUCTION BY LEVERAGING MCMC

We conduct dataset construction in two stages: 1) using a masked sentence as the prior and employing MCMC to find memorized image trigger prompts, and 2) using the memorized image trigger prompts as the prior for augmentation through MCMC.

Using Masked Sentence as Prior. This stage aims to discover new memorized images. The sentence is initialized with sentence of length n [MASK] token, *i.e.* $\mathbf{p}_0 = \{[\text{MASK}], [\text{MASK}], \dots, [\text{MASK}]\}$. We then employ Algorithm 1 initialized with \mathbf{p}_0 to obtain the candidate prompt. Similar to the conventional approach (Carlini et al., 2023) to extract train images, we then generate 100 images for this prompt and leverage DBSCAN (Ester et al., 1996) clustering algorithm with SSCD (Pizzi et al., 2022) to extract images forming at least 20 nodes. Those images are employed to Reverse Image Search API to find train image sources and human verification is conducted.

Using Found Trigger Prompts as Prior. This stage aims to augment memorized image trigger prompts. We leverage the prompts found in the previous stage or those provided by Webster (2023) as the prior, π_0 . In this process, we employ a slightly modified algorithm to enhance diversity. Instead of running a single chain for one prompt, we run n separate chains for each word position in an n -length sentence, treating each position as the first updating index in Gibbs sampling. This method ensures a varied exploration of the prompt space. We then save the top 100 prompts with the highest $D_\theta(\cdot)$. We retained all prompts generated during the MCMC sampling process and then selected 20 augmented prompts per original prompt, considering diversity. The detailed process is provided in Appendix C.

4 DATASET STATISTICS AND EFFICIENCY OF THE PROPOSED ALGORITHM

4.1 DATASET STATISTICS

Table 1 presents the number of memorized images and trigger prompts obtained using the methodology described in Section 3. For both Stable Diffusion 1 and 2, the number of prompts has increased more than fivefold for each model and more than 9 times in total compared to those reported by Webster (2023). The number of memorized images

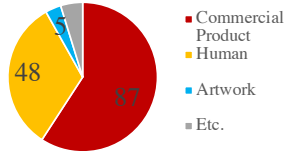


Figure 1: Components of Memorized Images in Stable Diffusion 1.

Table 1: Comparison of the number of memorized images and trigger prompts in each dataset. Our dataset is significantly larger in terms of the number of trigger prompts across all models. Please note that images sharing the same layout, as shown in Figure 3, have been counted as a single image.

	Stable Diffusion 1		Stable Diffusion 2		DeepFloydIF		Realistic Vision	
	Trigger Prompt #	Mem. Image #	Trigger Prompt #	Mem. Image #	Trigger Prompt #	Mem. Image #	Trigger Prompt #	Mem. Image #
Webster (2023)	325	111	210	25	162	17	354	119
MemBench	3000	151	1500	55	309	51	1352	148

Table 2: Comparison of the efficiency of our method and other prompt space optimization methods. Experiment was done on 1 A100 GPU. “-” denotes the failure of the valid search.

	Greedy Search	ZeroCap	PEZ	ConZIC	Ours
Hours/Memorized Image	5.7	-	-	3.81	2.08

included in the dataset has also increased, with Stable Diffusion 2 showing an increase of over twofold. Additionally, we provide memorized images and trigger prompts for DeepFloydIF (Shonenkov et al., 2023), which has a cascaded structure, and Realistic Vision (CivitAI, 2023), an open-source diffusion model. For these two models, we provide a larger number of memorized images and trigger prompts than Webster et al. We have also applied our algorithm to the more recent model, Stable Diffusion 3 (Esser et al., 2024). Please refer to Appendix E for the results. The composition of the images included in MemBench is shown in Figure 1, illustrating that the memorized images encompass a substantial number of commercial product images and human images. It also includes artwork such as brand logos.

4.2 EFFICACY OF MEMORIZED IMAGE TRIGGER PROMPT SEARCHING

In this section, we validate the efficiency of our method in discovering memorized images without access to \mathcal{D}_{train} . The task of finding memorized image trigger prompts without \mathcal{D}_{train} is defined as follows: without any prior information, the method must automatically find trigger prompts that induce memorized images. This involves: 1) selecting candidate prompts, 2) generating 100 images for each candidate prompt, 3) applying DBSCAN (Ester et al., 1996) clustering with SSCD to get candidate images and using a Reverse Image Search API to verify those images’ presence on the web. To the best of our knowledge, this task is novel, so we provide naive baselines. As the first baseline, we perform a greedy search by measuring D_θ for all prompts in the prompt dataset and selecting the top 200 prompts as candidate prompts. For the prompt dataset, we leveraged DiffusionDB, which contains 13M prompts collected from diffusion model users. Additionally, we provide three other baselines, all of which are algorithms that solve optimization problems in the prompt space. For two of these baselines, we adapt ZeroCap (Tewel et al., 2022) and ConZIC (Zeng et al., 2023), methods designed to maximize the CLIP Score for zero-shot image captioning, by replacing their objective function with D_θ . Similarly, we also adapt PEZ (Wen et al., 2024a) by substituting its objective function with D_θ to serve as another baseline. For each of these three methods, we conducted 200 iterations and obtained prompts. For our method, we performed 200 MCMC runs with 150 iterations each, and selected the resulting prompts as candidate prompts. For more detailed implementation, please refer to Appendix G.

The results are shown in Table 2. Our method significantly outperforms other methods. The results in Table 2 demonstrate that our method significantly outperforms others. To generate 200 candidate prompts, ZeroCap and PEZ required 44 and 33 hours, respectively, on an A100 GPU but failed to identify any memorized image trigger prompts. ZeroCap’s sequential prediction hindered the prompt being optimized to have higher D_θ values than general prompts. For PEZ, we observed that the prompts were optimized to produce images with a specific color (e.g., sunflower fields, grassy fields), and the prompts themselves were very unnatural. ConZIC identified 6 memorized images in 24 hours but struggled with local minima and lacked diversity in its optimization process, resulting in lower efficiency compared to our method.



Figure 2: The necessity of measuring the Aesthetic score. Images generated with the mitigation method applied are not desirable but achieve a low SSCD while maintaining a high CLIP Score.

Comparing with baselines (Carlini et al., 2023; Webster, 2023) that leverage LAION itself is challenging, as the dataset is no longer available and the elements for reimplementation are omitted in the corresponding papers. However, as mentioned in Section 2, their memory-inefficient and computationally intensive methods provided only a few memorized images and trigger prompts.

5 MEMBENCH: METRICS, SCENARIOS AND REFERENCE PERFORMANCE

Metrics. We present rigorous metrics for correctly evaluating mitigation methods, which include similarity score, Text-Image alignment score, and quality score. Following previous works, we adopt **SSCD** (Pizzi et al., 2022) as the similarity score and measure max SSCD between a generated image using trigger prompt and memorized images. In detail, if a prompt \mathbf{p}^* triggers images $\{\mathbf{x}_1^*, \dots, \mathbf{x}_k^*\}$ included in \mathcal{D}_{train} , we measure $\max_{\mathbf{x} \in \{\mathbf{x}_1^*, \dots, \mathbf{x}_k^*\}} SSCD(\mathbf{x}(\mathbf{p}^*, \epsilon_\theta), \mathbf{x})$. Secondly, we adopt **CLIP Score** (Hessel et al., 2021) to measure Text-Image alignment between prompt and generated images. Lastly, We adopt an **Aesthetic Score** (Schuhmann et al., 2022) as the image quality score. While previous works did not measure image quality scores, we observed issues shown in Figure 2. When memorization mitigation methods are applied, we observed that image quality degrades, the rich context generated by the diffusion model is destroyed, or distorted images are formed. **To further investigate, we have quantified this by calculating the standard deviation of Aesthetic Score.** An ideal memorization mitigation method should be able to preserve the generation capabilities of the diffusion model.

Scenarios. To ensure that memorization mitigation methods can be generally applied to diffusion models, we provide two scenarios: the memorized image trigger prompt scenario and the general prompt scenario. First, the memorized image trigger prompt scenario evaluates whether mitigation methods can effectively prevent the generation of memorized images. This scenario uses the memorized image trigger prompts we identified in Section 3. We generate 10 images for each trigger prompt and measure the Top-1 SSCD and the mean values of the Top-3 SSCD. We also measure the proportion of images with SSCD exceeding 0.5. For CLIP Score and Aesthetic Score, we calculate the average value across all generated images. Second, the general prompt scenario ensures that the performance of the diffusion model does not degrade when using prompts other than trigger prompts. We leverage the COCO (Lin et al., 2014) validation set as general prompts. In this scenario, images are generated once per prompt, and the average CLIP Score and Aesthetic Score are measured.

Reference Performance. We propose a reference performance for interpreting the performance of mitigation methods. An effective mitigation method should be able to reduce SSCD while maintaining CLIP Score. However, although SSCD is a metric designed to compare the structural similarity of images for copy detection tasks, it inevitably includes semantic meaning due to the self-supervised nature of the trained neural network. On the other hand, the semantic meaning of the trigger prompt should still be reflected in the generated image to maintain CLIP Score even when a mitigation method is applied. Therefore, it is uncertain how much the SSCD between memorized images and generated images can be reduced while maintaining the CLIP Score between trigger prompts and generated images. In this regard, we provide a reference performance to indicate how much SSCD can be reduced while maintaining a high CLIP Score. We assume querying images with trigger prompts via the Google Image API⁴ as a strong proxy model for the generative model and provide the reference performance based on this approach. Please refer to Appendix D for details.

⁴<https://developers.google.com/custom-search>

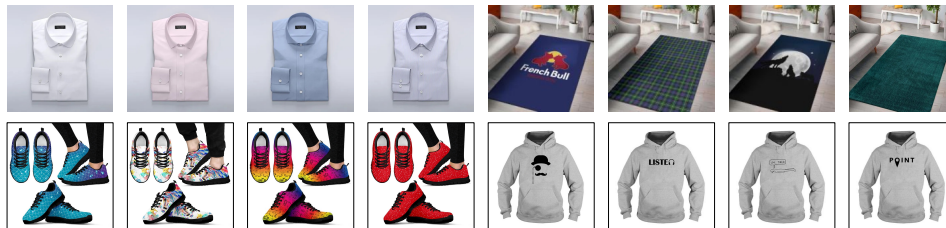


Figure 3: Results of images found by leveraging Reverse Image Search API to the images generated from trigger prompts. The shared layout suggests the occurrence of image memorization.

Table 3: Performance evaluation of image memorization mitigation methods in MemBench. Please refer to Appendix G.3 for the details of hyper-parameters.

		MemBench						COCO		
		Top-1 SSCD ↓	Top-3 SSCD ↓	SSCD > 0.5 ↓	CLIP ↑	Aesth. ↑	Aesth. std. ↓	CLIP ↑	Aesth. ↑	Aesth. std. ↓
Base		0.641	0.605	0.451	0.273	5.25	0.43	0.321	5.37	0.36
Reference Performance	(API search)	0.088	-	-	0.310	-	-	-	-	-
RNA(Somepalli et al., 2023b)	n = 1	0.479	0.425	0.241	0.270	5.18	0.53	0.314	5.34	0.36
	n = 2	0.389	0.338	0.165	0.270	5.14	0.55	0.310	5.33	0.37
	n = 3	0.329	0.280	0.121	0.267	5.13	0.56	0.307	5.30	0.37
	n = 4	0.287	0.239	0.089	0.264	5.10	0.58	0.304	5.29	0.39
	n = 5	0.254	0.213	0.074	0.262	5.08	0.59	0.302	5.28	0.39
	n = 6	0.228	0.189	0.055	0.258	5.06	0.59	0.298	5.24	0.38
RTA (Somepalli et al., 2023b)	n = 1	0.497	0.446	0.265	0.269	5.20	0.52	0.316	5.34	-
	n = 2	0.397	0.347	0.175	0.268	5.19	0.53	0.314	5.32	0.36
	n = 3	0.330	0.285	0.129	0.266	5.17	0.54	0.310	5.29	0.36
	n = 4	0.282	0.240	0.094	0.264	5.15	0.55	0.306	5.27	0.37
	n = 5	0.257	0.217	0.080	0.262	5.14	0.53	0.302	5.26	0.37
	n = 6	0.228	0.190	0.056	0.258	5.10	0.56	0.299	5.27	0.38
Wen et al. (2024b)	l = 7	0.410	0.346	0.134	0.270	5.16	0.54	0.321	5.37	0.36
	l = 6	0.355	0.289	0.089	0.270	5.15	0.55	0.321	5.37	0.36
	l = 5	0.312	0.246	0.059	0.269	5.14	0.56	0.321	5.37	0.36
	l = 4	0.259	0.199	0.035	0.268	5.13	0.57	0.321	5.37	0.36
	l = 3	0.181	0.139	0.015	0.264	5.11	0.59	0.321	5.37	0.36
	l = 2	0.096	0.075	0.001	0.242	4.97	0.64	0.321	5.37	0.36
Ren et al. (2024)	c = 1.0	0.289	0.247	0.083	0.263	5.17	0.57	0.316	5.33	0.38
	c = 1.1	0.283	0.239	0.071	0.260	5.17	0.57	0.313	5.31	0.38
	c = 1.2	0.278	0.232	0.058	0.257	5.15	0.58	0.309	5.28	0.39
	c = 1.3	0.275	0.227	0.050	0.254	5.14	0.58	0.304	5.26	0.39

6 DEEPER ANALYSIS INTO IMAGE MEMORIZATION

Secondly, we explore the cause of image memorization in Stable Diffusion 2, trained on LAION-5B, whose duplicates are removed. Previous works (Somepalli et al., 2023b; Gu et al., 2023) suggested that image memorization issues arise from duplicate images in the training data. Webster et al. (2023) confirmed that the LAION-2B dataset contains many duplicate images likely to be memorized. However, Stable Diffusion 2 still exhibits image memorization issues while reduced. We hypothesize that this memorization arises due to layout duplication. Figure 3 shows the images found by Reverse Image Search API that are memorized by Stable Diffusion. We found that there are often over 100 images on the web with the same layout but different color structures. LAION-5B underwent deduplication based on URLs⁵, but this process may not have removed these images. These layout memorizations are also obviously subject to copyright, posing potential social issues. Additional examples are provided in Appendix H.

7 EVALUATION OF IMAGE MEMORIZATION MITIGATION METHODS

In this section, we evaluate image memorization mitigation methods on our MemBench in Stable Diffusion 1. For results of Stable Diffusion 2, please refer to Appendix F.2.

Baselines. We use Stable Diffusion 1.4 as the base model. The image memorization mitigation methods evaluated include: 1) RTA (Somepalli et al., 2023b), which applies random token insertion to the prompt, 2) RNA (Somepalli et al., 2023b), which inserts a random number between $[0, 10^6]$ into the prompt, 3) method proposed by Wen et al. (2024b) that applies adversarial attacks to text

⁵<https://laion.ai/blog/laion-5b/>

486 embeddings, and 4) method proposed by Ren et al. (2024) that rescales cross-attention. Image
487 generation is performed using the DDIM (Song et al., 2021a) Scheduler with a guidance scale of 7.5
488 and 50 inference steps.

489
490 **Results.** We present the experimental results in Table 3. As shown in Table 3, for all methods,
491 lowering the SSCD significantly reduces both the CLIP Score and the Aesthetic Score. This indicates
492 a degradation in text-image alignment and image quality. In particular, upon examining images with
493 low Aesthetic Scores, we observe that issues in Figure 2 occur across all methods. While Ren et al.
494 (2024) measured FID, they reported that FID decreases when their method is applied. They attribute
495 this phenomenon to the mitigation method preventing memorized images from being generated,
496 thereby increasing the diversity of generated images. As a result, FID does not effectively measure
497 image quality. We provide FID values in Appendix F.1. However, the Aesthetic Score offers a more
498 straightforward way to evaluate individual image quality and better highlight image quality issues.
499 Moreover, when hyper-parameters are set as high values for mitigation methods, it leads not only to
500 a lower Aesthetic Score but also to a much larger standard deviation. This indicates that diffusion
501 model outputs become unreliable. As reported by Wen et al. (2024b), all methods exhibit a trade-off
502 between SSCD and CLIP Score. Regarding the reference performance obtained via API search, it
503 can be observed that the SSCD can be reduced to 0.088 while maintaining a high CLIP Score. Due to
504 the inherent limitations of the Stable Diffusion baseline model, the CLIP Score cannot exceed 0.273
505 when mitigation methods are applied. However, mitigation methods should aim to reduce the Top-1
506 SSCD to around 0.088 while maintaining at least this level of CLIP Score.

506 To provide a more detailed analysis of each method, we observe that the approach proposed by Wen et
507 al. achieves the best performance in the trade-off between SSCD and CLIP Score. However, to
508 reduce the proportion of images with SSCD exceeding 0.5—indicative of image memorization—to
509 nearly zero, their method still requires a reduction in CLIP Score by 0.025. Given the scale of the
510 CLIP Score, this drop suggests that the generated images may be only marginally related to the given
511 prompts. Moreover, a significant decrease in the Aesthetic Score is also observed. On the other hand,
512 the method proposed by Wen et al. has an additional advantage: it does not result in any performance
513 drop in the general prompt scenario on the COCO dataset, making it the most suitable option for
514 practical applications as of now.

515 The most recent method proposed by Ren et al. (2024) shows a considerable reduction in the CLIP
516 Score. Even at the lowest hyper-parameter setting ($c = 1.0$), the reductions in both CLIP Score and
517 Aesthetic Score are substantial, limiting its general applicability to diffusion models. The most basic
518 approaches, RNA and RTA, show a decrease in CLIP Score by 0.015 at the hyper-parameter setting
519 ($n = 6$) that lowers the proportion of images with SSCD exceeding 0.5 to 0.05. This is expected,
520 given the nature of these methods: both attempt to prevent image memorization by adding irrelevant
521 tokens to the prompts. As a result, RNA and RTA are unreliable for application to diffusion models.

522 523 8 CONCLUSION 524 525

526 We have presented MemBench, the first benchmark for evaluating image memorization mitigation
527 methods in diffusion models. MemBench includes various memorized image trigger prompts,
528 appropriate metrics, and a practical scenario to ensure that mitigation methods can be effectively
529 applied in practice. We have provided the reference performance that mitigation methods should aim
530 to achieve. Through MemBench, we have confirmed that existing image memorization mitigation
531 methods are still insufficient for application to diffusion models in practical scenarios. The lack of a
532 benchmark may have previously hindered the research of effective mitigation methods. However, we
533 believe that our benchmark will facilitate significant advancements in this field.

534 **Limitations and Future Work.** Another contribution of our work is providing an algorithm for
535 efficiently searching memorized image trigger prompts based on MCMC. Our approach is faster
536 than other searching algorithms we have tried, yet it does not exhibit exceptionally high speed.
537 Consequently, due to time constraints, we were unable to provide a larger number of memorized
538 images. However, our method allows for the continuous search of more memorized images and their
539 corresponding trigger prompts, and we plan to update the dataset regularly. Additionally, we aim to
enhance the efficiency of our memorized image trigger prompt searching algorithm in the future.

540 ETHICS STATEMENT

541

542 Our work introduces a technique for extracting the training data of diffusion models. This could
 543 potentially harm the rights of model owners or image copyright holders. Therefore, it is crucial to
 544 handle this technique with caution to avoid any infringement issues. For more details, please refer to
 545 Appendix I.

546

547 REPRODUCIBILITY STATEMENT

548

549 We provide the code for our training data extraction algorithm, the dataset, and the evaluation in the
 550 supplementary material.

551

552 REFERENCES

553

554 Nicolas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Sehwal, Florian Tramer, Borja
 555 Balle, Daphne Ippolito, and Eric Wallace. Extracting training data from diffusion models. In *32nd*
 556 *USENIX Security Symposium (USENIX Security 23)*, 2023.

557

558 Yunhao Chen, Xingjun Ma, Difan Zou, and Yu-Gang Jiang. Extracting training data from uncondi-
 559 tional diffusion models. *arXiv preprint arXiv:2406.12752*, 2024.

560

561 CivitAI. CivitAI, 2023.

562

563 Giannis Daras, Kulin Shah, Yuval Dagan, Aravind Gollakota, Alex Dimakis, and Adam Klivans.
 564 Ambient diffusion: Learning clean distributions from corrupted data. In *Advances in Neural*
Information Processing Systems (NeurIPS), 2024.

565

566 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep
 567 bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.

568

569 Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam
 570 Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion En-
 571 glish, Kyle Lacey, Alex Goodwin, Yannik Marek, and Robin Rombach. Scaling rectified flow
 572 transformers for high-resolution image synthesis. *arXiv preprint arXiv:2403.03206*, 2024.

573

574 Martin Ester, Hans-Peter Kriegel, Jörg Sander, Xiaowei Xu, et al. A density-based algorithm for
 575 discovering clusters in large spatial databases with noise. In *Proceedings of the 2nd International*
Conference on Knowledge Discovery and Data Mining (KDD), 1996.

576

577 Stuart Geman and Donald Geman. Stochastic relaxation, gibbs distributions, and the bayesian
 578 restoration of images. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*,
 579 (6):721–741, 1984.

580

581 Xiangming Gu, Chao Du, Tianyu Pang, Chongxuan Li, Min Lin, and Ye Wang. On memorization in
 582 diffusion models. *arXiv preprint arXiv:2310.02664*, 2023.

583

584 W. K. Hastings. Monte carlo sampling methods using markov chains and their applications.
 585 *Biometrika*, 57(1):97–109, 1970.

586

587 Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. CLIPScore: a reference-
 588 free evaluation metric for image captioning. In *Conference on Empirical Methods in Natural*
Language Processing (EMNLP), 2021.

589

590 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*,
 591 2022.

592

593 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in*
Neural Information Processing Systems (NeurIPS), 2020.

594

595 Alicia A. Johnson, Galin L. Jones, and Ronald C. Neath. Component-wise markov chain monte
 596 carlo: Uniform and geometric ergodicity under mixing and composition. *Statistical Science*, 28(3):
 597 360–375, 2013.

- 594 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr
595 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–*
596 *ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings,*
597 *Part V 13*, pp. 740–755. Springer, 2014.
- 598
599 Xiao Liu, Xiaoliu Guan, Yu Wu, and Jiaxu Miao. Iterative ensemble training with anti-gradient
600 control for mitigating memorization in diffusion models. *arXiv preprint arXiv:2407.15328*, 2024.
- 601 Omiros Papaspiliopoulos and Gareth O. Roberts. Stability of the gibbs sampler for bayesian hierar-
602 chical models. *The Annals of Statistics*, 36(1):95–117, 2008.
- 603
604 Ed Pizzi, Sreya Dutta Roy, Sugosh Nagavara Ravindra, Priya Goyal, and Matthijs Douze. A self-
605 supervised descriptor for image copy detection. In *Proceedings of the IEEE/CVF Conference on*
606 *Computer Vision and Pattern Recognition*, pp. 14532–14542, 2022.
- 607 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
608 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
609 models from natural language supervision. In *International Conference on Machine Learning*
610 *(ICML)*, 2021.
- 611
612 Jie Ren, Yaxin Li, Shenglai Zen, Han Xu, Lingjuan Lyu, Yue Xing, and Jiliang Tang. Unveiling and
613 mitigating memorization in text-to-image diffusion models through cross attention. *arXiv preprint*
614 *arXiv:2403.11052*, 2024.
- 615 Christian P Robert, George Casella, and George Casella. *Monte Carlo statistical methods*, volume 2.
616 Springer, 1999.
- 617
618 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
619 resolution image synthesis with latent diffusion models. In *IEEE Conference on Computer Vision*
620 *and Pattern Recognition (CVPR)*, 2022.
- 621 Gustavo Santana. Gustavosta/stable-diffusion-prompts · datasets at hugging face, De-
622 cember 2022. URL [https://huggingface.co/datasets/Gustavosta/](https://huggingface.co/datasets/Gustavosta/Stable-Diffusion-Prompts)
623 [Stable-Diffusion-Prompts](https://huggingface.co/datasets/Gustavosta/Stable-Diffusion-Prompts).
- 624
625 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi
626 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An
627 open large-scale dataset for training next generation image-text models. 2022.
- 628
629 Alex Shonenkov et al. Deep image floyd. <https://github.com/deep-floyd/IF>, 2023.
- 630 Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Diffusion
631 art or digital forgery? investigating data replication in diffusion models. In *IEEE Conference on*
632 *Computer Vision and Pattern Recognition (CVPR)*, 2023a.
- 633
634 Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Understand-
635 ing and mitigating copying in diffusion models. In *Advances in Neural Information Processing*
636 *Systems (NeurIPS)*, 2023b.
- 637
638 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *Interna-*
639 *tional Conference on Learning Representations (ICLR)*, 2021a.
- 640 Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben
641 Poole. Score-based generative modeling through stochastic differential equations. In *International*
642 *Conference on Learning Representations (ICLR)*, 2021b.
- 643
644 Alexander Terenin, Daniel Simpson, and David Draper. Asynchronous gibbs sampling. In *Interna-*
645 *tional Conference on Artificial Intelligence and Statistics*, pp. 144–154. PMLR, 2020.
- 646
647 Yoad Tewel, Yoav Shalev, Idan Schwartz, and Lior Wolf. Zerocap: Zero-shot image-to-text generation
for visual-semantic arithmetic. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
and Pattern Recognition, pp. 17918–17928, 2022.

648 Ryan Webster. A reproducible extraction of training images from diffusion models. *arXiv preprint*
649 *arXiv:2305.08694*, 2023.

650
651 Ryan Webster, Julien Rabin, Loic Simon, and Frederic Jurie. On the de-duplication of laion-2b. *arXiv*
652 *preprint arXiv:2303.12733*, 2023.

653 Yuxin Wen, Neel Jain, John Kirchenbauer, Micah Goldblum, Jonas Geiping, and Tom Goldstein.
654 Hard prompts made easy: Gradient-based discrete optimization for prompt tuning and discovery.
655 *Advances in Neural Information Processing Systems*, 36, 2024a.

656
657 Yuxin Wen, Yuchen Liu, Chen Chen, and Lingjuan Lyu. Detecting, explaining, and mitigating
658 memorization in diffusion models. In *International Conference on Learning Representations*
659 *(ICLR)*, 2024b.

660 Zequn Zeng, Hao Zhang, Ruiying Lu, Dongsheng Wang, Bo Chen, and Zhengjue Wang. Conzic:
661 Controllable zero-shot image captioning by sampling-based polishing. In *Proceedings of the*
662 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 23465–23476, 2023.

663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

Appendix

702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

CONTENTS

1	Introduction	1
2	Related Work	2
3	Searching Memorized Image Trigger Prompt with MCMC	3
3.1	Preliminary	3
3.2	Memorization Trigger Prompt Searching as an Optimization Problem	4
3.3	Constructing MCMC by Leveraging D_θ	5
3.4	Dataset Construction by Leveraging MCMC	6
4	Dataset Statistics and Efficiency of the Proposed Algorithm	6
4.1	Dataset Statistics	6
4.2	Efficacy of Memorized Image Trigger Prompt Searching	7
5	MemBench: Metrics, Scenarios and Reference Performance	8
6	Deeper Analysis into Image Memorization	9
7	Evaluation of Image Memorization Mitigation Methods	9
8	Conclusion	10
A	Detailed Derivation of Transition Matrix	16
B	Stable Diffusion replicating commercial products currently on sale	16
C	Data Construction Details	17
C.1	Data Augmentation Leveraging MCMC	17
C.2	Data Augmentation Performance	17
D	Reference Performance Based on Image Search API	19
E	Extension to Stable Diffusion 3	19
F	Evaluation of Image Memorization Mitigation Method on MemBench	20
F.1	FID of mitigation methods measured on MemBench	20
F.2	Evaluation of Image Memorization Mitigation Method on Stable Diffusion 2	20
G	Details of Experiments	21
G.1	Implementation Details of Baselines in Memorized Image Trigger Prompt Searching Experiment	21

756	G.2 Hyper-Parameters in Image Memorization Mitigation Methods	22
757	G.3 Hyper-Parameters in Memorized Image Trigger Prompt Searching Leveraging MCMC	22
758		
759		
760	H Additional Examples of Memorized Images	22
761		
762	I Datasheet	26
763		
764	I.1 Motivation	26
765	I.2 Composition	26
766	I.3 Collection Process	28
767		
768	I.4 Preprocessing, Cleaning, and/or Labeling	29
769	I.5 Uses	29
770		
771	I.6 Distribution and License	30
772	I.7 Maintenance	30
773		
774		
775		
776		
777		
778		
779		
780		
781		
782		
783		
784		
785		
786		
787		
788		
789		
790		
791		
792		
793		
794		
795		
796		
797		
798		
799		
800		
801		
802		
803		
804		
805		
806		
807		
808		
809		

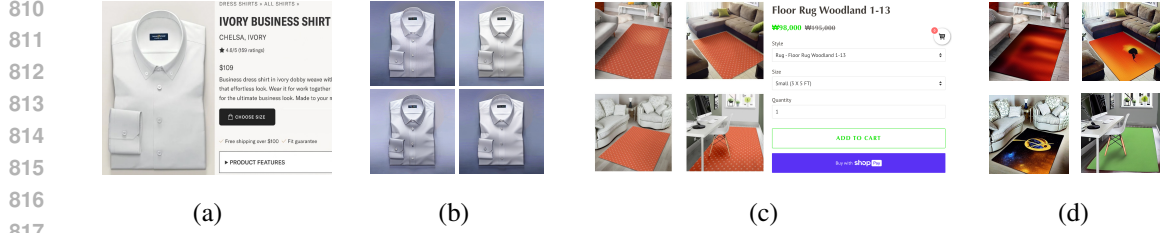


Figure 4: Examples of memorized images found using the Reverse Image Search API. (a), (c) Shirt/rug currently sold commercially, (b), (d) four images generated by Stable Diffusion

A DETAILED DERIVATION OF TRANSITION MATRIX

In this section, we provide a detailed derivation of the transition matrix that was omitted in Section 3. To recap, the transition matrix of the random scan Gibbs sampler for sampling the target distribution π is defined as follows:

$$\mathbf{T} = \sum_{i=1}^n \frac{1}{n} \cdot \mathbf{T}_i, \quad (10)$$

$$[\mathbf{T}_i]_{\mathbf{p}^j \rightarrow \mathbf{p}^{j+1}} = \begin{cases} \pi(\mathbf{p}_i^{j+1} | \mathbf{p}_{-i}^j) & \text{if } \mathbf{p}_{-i}^j = \mathbf{p}_{-i}^{j+1} \\ 0 & \text{else,} \end{cases} \quad (11)$$

where n is the total length of sentence, $\mathbf{p}_{-i} = \{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_{i-1}, \mathbf{p}_{i+1}, \dots, \mathbf{p}_n\}$ and \mathbf{p}^j is a j -th state prompt. We proceed to derive the conditional probability distribution of the target distribution $\pi(\mathbf{p}) = \frac{e^{D_{\theta}(\mathbf{p})/K}}$:

$$\pi(\mathbf{p}_i^{j+1} | \mathbf{p}_{-i}^j) = \frac{\pi(\mathcal{P}_i = \mathbf{p}_i^{j+1}, \mathcal{P}_{-i} = \mathbf{p}_{-i}^j)}{\pi(\mathbf{p}_{-i}^j)} \quad (12)$$

$$= \frac{\pi(\mathcal{P}_i = \mathbf{p}_i^{j+1}, \mathcal{P}_{-i} = \mathbf{p}_{-i}^j)}{\sum_{\mathbf{w} \in \mathcal{W}} \pi(\mathcal{P}_i = \mathbf{w}, \mathcal{P}_{-i} = \mathbf{p}_{-i}^j)} \quad (13)$$

$$= \frac{e^{D_{\theta}(\mathcal{P}_i = \mathbf{p}_i^{j+1}, \mathcal{P}_{-i} = \mathbf{p}_{-i}^j)/K}}{\sum_{\mathbf{w} \in \mathcal{W}} e^{D_{\theta}(\mathcal{P}_i = \mathbf{w}, \mathcal{P}_{-i} = \mathbf{p}_{-i}^j)/K}} \quad (14)$$

By substituting Equation 14 into Equation 11, we ultimately derive the transition matrix as defined earlier in Equation 9.

$$[\mathbf{T}]_{\mathbf{p}^j \rightarrow \mathbf{p}^{j+1}} = \begin{cases} \frac{1}{n} \cdot \left(\frac{e^{D_{\theta}(\mathcal{P}_i = \mathbf{p}_i^{j+1}, \mathcal{P}_{-i} = \mathbf{p}_{-i}^j)/K}}{\sum_{\mathbf{w} \in \mathcal{W}} e^{D_{\theta}(\mathcal{P}_i = \mathbf{w}, \mathcal{P}_{-i} = \mathbf{p}_{-i}^j)/K}} \right) & \text{if } \mathbf{p}_{-i}^j = \mathbf{p}_{-i}^{j+1}, \\ 0 & \text{else.} \end{cases} \quad (15)$$

B STABLE DIFFUSION REPLICATING COMMERCIAL PRODUCTS CURRENTLY ON SALE

In this section, we provide a deeper analysis of the memorized images and trigger prompts in MemBench. We have found that **Stable Diffusion regenerates commercial products currently on sale**. While the possibility that diffusion models could memorize commercial images has been suggested (Carlini et al., 2023; Somepalli et al., 2023a), we are the first to confirm this. Unlike the previous studies (Carlini et al., 2023; Webster, 2023) that used image retrieval from LAION to find memorized images, we leverage a Reverse Image Search API to find those, which enable us this

864 verification. As shown in Figure 4.b, Stable Diffusion replicates images of commercially available
 865 shirts when given a specific prompt. Figure 4.d further illustrates the replication of layouts; for a
 866 commercially sold carpet, all layouts have been reproduced.
 867

868 C DATA CONSTRUCTION DETAILS

869 In this section, we provide a detailed explanation of the data construction process described in
 870 Section 3.4. We explain 1) how memorized image trigger prompts and corresponding memorized
 871 images for Stable Diffusion 1 and 2 in Table 1 were found, 2) implementation details of the data
 872 augmentation algorithm using MCMC, and 3) its efficiency.
 873

874 For Stable Diffusion 1 and Realistic Vision, we initialized sentences with n -length mask tokens
 875 and implemented Algorithm 1 to find new memorized image trigger prompts and corresponding
 876 memorized images (please refer to Section 3.4 “Using Masked Sentence as Prior”). We then
 877 perform the MCMC process with \mathbf{p}_0 initialized by trigger prompts to perform the augmentation
 878 (please refer to Section 3.4 “Using Found Trigger Prompts as Prior”). For Stable Diffusion
 879 2 and DeepFloydIF, the process of finding trigger prompts using masked sentences was omitted.
 880 This was due to two reasons: firstly, the prediction accuracy of D_θ for memorized image trigger
 881 prompts is lower for these models. Secondly, as Stable Diffusion 2 is trained on the deduplicated
 882 LAION-5B and LAION-A, the memorized image trigger prompts are sparser, making optimization
 883 from a masked sentence initialization difficult. Therefore, for Stable Diffusion 2 and DeepFloydIF,
 884 only the trigger prompt augmentation algorithm was leveraged. The prompts were initialized in two
 885 ways before undergoing the data augmentation process: 1) using trigger prompts found from Stable
 886 Diffusion 1, and 2) using trigger prompts provided by Webster (2023). We further elaborate on the
 887 data augmentation process below.
 888

889 C.1 DATA AUGMENTATION LEVERAGING MCMC

890 Trigger prompt augmentation was carried out using a different approach from trigger prompt searching
 891 (Algorithm 1). The process of generating candidate trigger prompts through prompt augmentation is
 892 detailed in Algorithm 2. We initialized \mathbf{p}_0 with the trigger prompt itself and then performed MCMC.
 893 As explained in Section 3.4, in trigger prompt augmentation, we run n separate chains for each word
 894 position in an n -length sentence, treating each position as the first updating index in Gibbs sampling.
 895 During the MCMC process, all prompts with calculated D_θ values were stored in the prompt bank.
 896 Additionally, we adopted an early stop counter. The prompts returned by Algorithm 2 tend to have
 897 low diversity due to the nature of Gibbs Sampling. Therefore, Algorithm 3 is applied to all returned
 898 prompts to create a smaller, more diverse subset of prompts. Afterward, these prompts undergo an
 899 image generation process, followed by human verification, before being added to the dataset.
 900

901 C.2 DATA AUGMENTATION PERFORMANCE

902 We present an evaluation of Algorithm 2, our proposed
 903 method for augmenting memorized image trigger prompts.
 904 To assess the effectiveness of Algorithm 2, we examine
 905 whether the prompts generated during the algorithm’s ex-
 906 ecution indeed trigger memorized images. Although Al-
 907 gorithm 2 is designed to return only the top T candidate
 908 trigger prompts, for this experiment, we investigate all the
 909 prompts generated during the execution of Algorithm 2
 910 to measure its performance. Given the extensive time re-
 911 quired to verify all candidate trigger prompts, we present
 912 a toy experiment focusing on a specific prompt: “The no
 913 limits business woman podcast,” which generates an im-
 914 age identical to Figure 5.a. For the experiment, we set
 915 the hyperparameters of Algorithm 2 as follows: $K = 1.5$,
 916 $N = 50$, $Q = 200$, $\kappa = 3$, and $s = 3$. The experiment
 917 was conducted using a single A100 GPU.



Figure 5: Memorized image utilized for toy experiment. Each image refers to (a) train data image in Stable Diffusion, (b) generated image using Stable Diffusion. The SSCD between (a) and (b) is measured to be 0.707.

Algorithm 2 Memorized Image Trigger Prompt Augmentation via Gibbs Sampling

```

918 1: Input: Diffusion model  $\theta$ , BERT model  $\phi$ , initial sentence  $\mathbf{p}^0$  with length  $n$ , iteration number
919  $N$ , number of proposal words  $Q$ , termination threshold  $\kappa$ , early stop counter threshold  $s$ ,
920 hyperparameter  $K$ .
921 2: Initialize early stop counter  $c \leftarrow 0$ 
922 3: Initialize prompt bank  $\mathcal{B} \leftarrow \{\mathbf{p}^0\}$ 
923 4: for  $k = 0$  to  $n$  do
924 5:   for  $j = 0$  to  $N$  do
925 6:     if  $j = 0$  then
926 7:        $i \leftarrow k$ 
927 8:     else
928 9:       Randomly select index  $i \in \{1, \dots, n\}$ 
929 10:    end if
930 11:     $\mathcal{W}_Q \leftarrow \arg \text{top}_Q p_\phi(\mathbf{w} \mid \mathbf{p}_{-i}^j)$ 
931 12:     $p(\mathbf{p}_i^{j+1} \mid \mathbf{p}_{-i}^j) \leftarrow \frac{e^{D_\theta(\mathcal{P}_i=\mathbf{p}_i^{j+1}, \mathcal{P}_{-i}=\mathbf{p}_{-i}^j)/K}}}{\sum_{\mathbf{w} \in \mathcal{W}_Q} e^{D_\theta(\mathcal{P}_i=\mathbf{w}, \mathcal{P}_{-i}=\mathbf{p}_{-i}^j)/K}}$ 
932 13:     $\mathbf{p}_i^{j+1} \leftarrow \text{Sample from } p(\mathbf{p}_i^{j+1} \mid \mathbf{p}_{-i}^j)$ 
933 14:     $\mathbf{p}^{j+1} \leftarrow (\mathbf{p}_1^j, \mathbf{p}_2^j, \dots, \mathbf{p}_i^{j+1}, \dots, \mathbf{p}_n^j)$ 
934 15:    Add  $\{(\mathbf{p}_1^j, \mathbf{p}_2^j, \dots, \mathbf{p}_{i-1}^j, \mathbf{w}, \mathbf{p}_{i+1}^j, \dots, \mathbf{p}_n^j) \mid \forall \mathbf{w} \in \mathcal{W}_Q\}$  to  $\mathcal{B}$ 
935 16:    if  $D_\theta(\mathbf{p}^{j+1}) < \kappa$  then
936 17:       $c \leftarrow c + 1$ 
937 18:    else
938 19:       $c \leftarrow 0$ 
939 20:    end if
940 21:    if  $c > s$  then
941 22:      break
942 23:    end if
943 24:  end for
944 25: end for
945 26: return  $\mathcal{B}$ 

```

Algorithm 3 Diversity Sampling

```

949 1: Input: Text encoder  $\phi$ , augmented prompts  $\mathcal{B}$ , return prompts number  $N$ 
950 2: Randomly select  $\mathbf{p}^* \in \mathcal{B}$ 
951 3: Initialize return prompt list  $\mathcal{R} \leftarrow \{\mathbf{p}^*\}$ 
952 4: while  $|\mathcal{R}| < N$  do
953 5:    $\mathbf{p}^* \leftarrow \arg \min_{\mathbf{p} \in \mathcal{B}} \max_{\mathbf{p}_r \in \mathcal{R}} \frac{\phi(\mathbf{p}) \cdot \phi(\mathbf{p}_r)}{\|\phi(\mathbf{p})\| \|\phi(\mathbf{p}_r)\|}$ 
954 6:    $\mathcal{B} \leftarrow \mathcal{B} \setminus \{\mathbf{p}^*\}$ 
955 7:    $\mathcal{R} \leftarrow \mathcal{R} \cup \{\mathbf{p}^*\}$ 
956 8: end while
957 9: return  $\mathcal{R}$ 

```

For those prompts generated during Algorithm 2, we filtered only prompts that show $D_\theta(\mathbf{p}) > 5$ and generated 10 images for each. Then we measure the Top-1 SSCD (Pizzi et al., 2022) with the image in Figure 5.a. We found that there were 4217 unique prompts with a Top-1 SSCD exceeding 0.7, indicating that they replicate train data image (as seen in Figure 5). Algorithm 2 took 7 minutes on an A100 GPU, producing 4217 augmented trigger prompts within this time frame. In addition, we categorized these prompts based on the number of words changed from the original prompt. Specifically, there were 753 trigger prompts with one word changed, 1923 with two words changed, 1352 with three words changed, 179 with four words changed, and 10 with five words changed. Interestingly, even with the modification of five out of the six words in the sentence, the altered prompts can still effectively induce memorized images. This demonstrates our method’s efficiency in generating a large number of augmented trigger prompts in a short period.



984 Figure 6: Trigger prompt searched by our MCMC algorithm and generated images with corresponding
985 prompt. The repeated and very similar images strongly suggest the occurrence of memorization.
986 Notably, the last images show a tendency to repeat specific text at the bottom left. This is a common
987 feature of memorized images generated by diffusion models, where text, URLs, or similar elements
988 from the source image are replicated. This strongly suggests that these repeated images are indeed
989 memorized images.

991 D REFERENCE PERFORMANCE BASED ON IMAGE SEARCH API

992
993
994 In Section 5, we posited that the Google Image Search API⁶ serves as a strong proxy model for
995 the generative model and presented the measured reference performance. We provide the details
996 of this approach in this section. Before explaining our use of the Google Image Search API, it
997 is important to reiterate our goal in presenting reference performance: to determine the extent to
998 which the SSCD (Pizzi et al., 2022) between the memorized image and the generated image can be
999 minimized while maintaining the semantic content of the trigger prompt. To address this question, we
1000 utilized the Google Image Search API to measure reference performance as follows: 1) Query 100
1001 images using the memorized image trigger prompt via the API. 2) Measure the CLIP Score (Hessel
1002 et al., 2021) between the 100 images and the trigger prompt. 3) Retain only the image with the
1003 Top-1 CLIP Score. 4) Measure the SSCD between this retained image and the memorized image
1004 triggered by the prompt in Stable Diffusion. 5) Repeat steps 1-4 for all memorized image trigger
1005 prompts in MemBench. After completing these steps, we reported the average Top-1 CLIP Score
1006 and the average SSCD of images with the Top-1 CLIP Score in Section 7. Our findings show that
1007 the SSCD can be reduced to 0.200 while maintaining a CLIP Score of 0.329. This indicates that the
1008 minimum achievable SSCD with maintaining CLIP Score is 0.210. Therefore, we should strive to
1009 develop mitigation methods that achieve this or better. Please note that we did not measure Aesthetic
1010 Score (Schuhmann et al., 2022) and evaluate the reference performance in COCO (Lin et al., 2014)
1011 settings, since comparing them with the mitigation method is not meaningful.

1012 E EXTENSION TO STABLE DIFFUSION 3

1013
1014 We applied our algorithm to Stable Diffusion 3. However, as the training data for Stable Diffusion
1015 3 is publicly unknown (no information is available), we were unable to perform the verification
1016 process, Reverse Image Search API. Thus, the searched images and prompts for Stable Diffusion 3
1017 cannot serve as a memorization benchmark. To be used as a memorization benchmark, two critical
1018 steps are essential: 1) Candidate Trigger Prompt Search Step and 2) Verification Step, where we
1019 confirm whether the repeated images actually exist in the training data, thereby verifying that they
1020 are indeed memorized images, as noted in Section 3.2. Without the verification step, we are not sure
1021 whether the searched images and prompts are memorized. Nevertheless, in Figure 6, we present the
1022 trigger prompts and duplicated images identified by our MCMC algorithm for Stable Diffusion 3.
1023 Although we cannot verify them due to the aforementioned limitations, we believe they represent
1024 strong candidates for memorized images.

1025 ⁶<https://developers.google.com/custom-search>

Table 4: FID scores of mitigation methods measured on MemBench.

	Base	RTA (Somepalli et al., 2023b)	RNA (Somepalli et al., 2023b)	Wen et al. (2024b)	Ren et al. (2024)
FID ↓	116.07	75.33	85.32	64.21	86.79

Table 5: Performance evaluation of image memorization mitigation methods in MemoBench for Stable Diffusion 2.

		MemBench					COCO	
		Top-1 SSCD ↓	Top-3 SSCD ↓	SSCD > 0.5 ↓	CLIP ↑	Aesthetic ↑	CLIP ↑	Aesthetic ↑
Base		0.629	0.593	0.448	0.281	5.41	0.333	5.35
Reference Performance	(API search)	0.207	-	-	0.301	-	-	-
RNA(Somepalli et al., 2023b)	n = 1	0.568	0.525	0.349	0.278	5.34	0.328	5.35
	n = 2	0.539	0.491	0.289	0.276	5.30	0.326	5.35
	n = 3	0.501	0.446	0.224	0.273	5.24	0.324	5.35
	n = 4	0.453	0.395	0.161	0.271	5.21	0.322	5.35
	n = 5	0.424	0.368	0.130	0.270	5.19	0.320	5.34
RTA (Somepalli et al., 2023b)	n = 1	0.590	0.549	0.365	0.276	5.35	0.332	5.34
	n = 2	0.562	0.515	0.317	0.275	5.31	0.330	5.31
	n = 3	0.529	0.475	0.261	0.272	5.26	0.329	5.30
	n = 4	0.479	0.428	0.211	0.272	5.21	0.325	5.27
	n = 5	0.452	0.393	0.167	0.271	5.17	0.324	5.25
Wen et al. (2024b)	l = 70	0.577	0.535	0.311	0.273	5.30	0.333	5.35
	l = 60	0.553	0.502	0.251	0.269	5.26	0.333	5.35
	l = 50	0.501	0.423	0.154	0.263	5.19	0.333	5.35
	l = 40	0.398	0.322	0.065	0.253	5.15	0.333	5.35
Ren et al. (2024)	c = 1.0	0.592	0.556	0.419	0.273	5.40	0.331	5.36
	c = 1.1	0.586	0.548	0.391	0.270	5.39	0.326	5.33
	c = 1.2	0.580	0.539	0.349	0.267	5.37	0.320	5.30
	c = 1.3	0.574	0.529	0.295	0.262	5.34	0.313	5.27

F EVALUATION OF IMAGE MEMORIZATION MITIGATION METHOD ON MEMBENCH

F.1 FID OF MITIGATION METHODS MEASURED ON MEMBENCH

In this section, we present the FID values measured on MemBench when applying mitigation methods to Stable Diffusion 1. As shown in Table 4, FID values increase when mitigation methods are applied. This aligns with the findings reported by Ren et al. (2024), as FID also captures diversity. Since the Stable Diffusion model generates identical images for trigger prompts, the generated images exhibit low diversity, leading to higher FID values. In contrast, when mitigation methods are applied, memorized images are not generated, resulting in increased diversity and consequently lower FID values. Therefore, FID does not effectively measure image quality but rather measures diversity. Image quality should instead be assessed using the Aesthetic Score we propose.

F.2 EVALUATION OF IMAGE MEMORIZATION MITIGATION METHOD ON STABLE DIFFUSION 2

In this section, we evaluate image memorization mitigation methods on our MemBench in Stable Diffusion 2.

We present the experimental results in Table 5. When each memorization mitigation method is applied, although SSCD (Pizzi et al., 2022) is reduced, there is a drop in both CLIP Score (Hessel et al., 2021) and Aesthetic Score (Schuhmann et al., 2022). Additionally, compared to the reference performance provided by the Google Image Search API, the performance of these methods is insufficient. The method proposed by Wen et al. (2024b) shows less capability in reducing SSCD while maintaining the CLIP Score compared to RNA (Somepalli et al., 2023b) and RTA (Somepalli et al., 2023b) in the memorized image trigger prompt scenario. However, in the practical scenario of the COCO validation set, its performance remains equivalent to the base Stable Diffusion 2.

1080 G DETAILS OF EXPERIMENTS

1081 G.1 IMPLEMENTATION DETAILS OF BASELINES IN MEMORIZED IMAGE TRIGGER PROMPT 1082 SEARCHING EXPERIMENT

1083 In this section, we provide implementation details of other baselines that we have tried to search
1084 the memorized image trigger prompts, presented in Section 4.2. All three algorithms (Wen et al.,
1085 2024a; Zeng et al., 2023; Tewel et al., 2022) that we provide are originally intended to solve various
1086 optimization problems in the text space, using specific objective functions. For our experiments to
1087 search for memorized image trigger prompts, we replaced each method’s objective function with D_θ .

1088 **ZeroCap (Tewel et al., 2022).** ZeroCap is an optimization method developed for zero-shot image
1089 captioning tasks. This method leverages a pre-trained CLIP (Radford et al., 2021) to measure the
1090 CLIP similarity between an image and the current caption, and manipulate the prompt to maximize
1091 this CLIP Score, searching the best caption that describes the image. ZeroCap predicts the next word
1092 using a large language model (LLM) and sequentially adds tokens to the prompt in a manner that
1093 maximizes CLIP similarity. Additionally, a Context Cache is introduced for gradient descent, where
1094 the Context Cache is a set of key-value pairs derived when the current prompt is embedded into the
1095 LLM. The optimization function is consistute of 1) CLIP similarity loss between the image and the
1096 prompt, and 2) the cross-entropy (CE) loss between the distribution of the predicted token of the
1097 original Context Cache and that of the updated Context Cache. ZeroCap performs gradient descent
1098 on the optimization function to update the Context Cache five times, after which the token predicted
1099 by this Context Cache is designated as the next token to continue the sentence. Furthermore, beam
1100 search is utilized in this process. In our experiments, we replaced the CLIP similarity loss with D_θ
1101 and implemented the algorithm accordingly. All hyper-parameters were set to match those in the
1102 original paper.

1103 To generate 200 candidate prompts using ZeroCap, it took approximately 44 hours on an A100 GPU,
1104 yet not a single memorized image trigger prompt was found. While D_θ values were higher compared
1105 to those of general prompts (captions from COCO validation set), they were still lower than the values
1106 for actual trigger prompts. This suggests an inherent issue with ZeroCap’s sequential prediction
1107 method.

1108 **PEZ (Wen et al., 2024a).** The PEZ algorithm is an optimization technique designed to find prompts
1109 that will induce a diffusion model to generate a specific desired image. The algorithm operates in two
1110 main steps for each iteration: 1) perform gradient descent on the prompt in the continuous space with
1111 respect to the diffusion model’s CLIP model, and 2) project the updated prompt back into the discrete
1112 space of the CLIP’s embedding space. In our adaptation of this algorithm, we utilized $D_\theta(\mathbf{p})$ as the
1113 objective function for calculating the gradient. All hyper-parameters were set to match those in the
1114 original paper.

1115 To generate 200 candidate prompts using PEZ, it took approximately 33 hours on an A100 GPU.
1116 However, similar to ZeroCap, not a single memorized image trigger prompt was found. Although
1117 the D_θ values were comparable to those of actual trigger prompts, memorized images were not
1118 discovered. Upon inspection, we observed that the prompts were optimized to produce images with a
1119 specific color (e.g., sunflower fields, grassy fields), and the prompts themselves were very unnatural.
1120 This suggests that the optimization process did not result in the desired memorized image trigger
1121 prompts.

1122 **ConZIC (Zeng et al., 2023).** ConZIC, like ZeroCap, is a technique designed to optimize the CLIP
1123 Score for zero-shot image captioning tasks. Similar to our approach, ConZIC selects a single word
1124 within the sentence, predicts the word using BERT, and then replace it with the word which shows
1125 the highest value of objective function. The objective function here is a sum of the CLIP similarity
1126 and the conditional probability distribution from BERT. In our experiments, we substituted the CLIP
1127 similarity with D_θ as the objective function.

1128 To generate 200 candidate prompts using ConZIC, it took approximately 24 hours on an A100 GPU.
1129 Unlike the other methods, ConZIC successfully identified 6 memorized images. However, ConZIC’s
1130 optimization process is designed to consistently update the prompt to maximize the objective function,
1131 which tends to result in getting stuck in local minima and the lack of diversity. These lead to less
1132 efficiency compared to our method.

Table 6: Hyper-parameters leveraged in memorized image trigger prompt searching using our algorithm. Here, n represents the sentence length, N is the iteration number, Q denotes the number of proposal words, K stands for the temperature, κ is the termination threshold, s is the early stop counter threshold, and T is the number of return candidate prompts.

Model	Method	n	N	Q	K	κ	s	T
Stable Diffusion 1	Algorithm 1	8	150	200	0.1	5	-	-
	Algorithm 2	-	20	200	1.5	3	3	100
Stable Diffusion 2	Algorithm 2	-	20	200	5.0	50	3	100

G.2 HYPER-PARAMETERS IN IMAGE MEMORIZATION MITIGATION METHODS

In Section 7, we evaluated the performance of image memorization mitigation methods on MemBench and presented the results in Table 3. However, due to space constraints, we omitted the explanations of various hyper-parameters in the table. Here, we provide a detailed explanation of these hyper-parameters. Firstly, RTA (Somepalli et al., 2023b) and RNA (Somepalli et al., 2023b) are methods that insert random words or numbers into the prompt. The parameter n in the table indicates the number of words or numbers inserted. The method proposed by Wen et al. (2024b) involves updating the prompt embedding to minimize D_θ . Here, the threshold for lowering $D_\theta(\mathbf{p})$, denoted as the early stopping loss l , becomes a hyper-parameter, *i.e.* the prompt \mathbf{p} is updated until $D_\theta(\mathbf{p}) < l$. All other hyper-parameters followed the settings in the original paper: an Adam optimizer with a learning rate of 0.05 and a maximum of 10 steps was used for training. Ren et al. (2024) provides a method that inversely amplifies the attention score for the beginning token by adjusting the input logits of the softmax operator in the cross-attention. To be precise, let the original input logits be denoted as $\mathbf{s} = (s_1, s_2, \dots, s_N)$, where s_i is the logit of the i -th token. The re-scaled logit vector \mathbf{s}' is:

$$\mathbf{s}' = (Cs_1, s_2, \dots, s_{N-S}, -\infty, \dots, -\infty). \quad (16)$$

Here, the scale factor C for the beginning token s_1 becomes a hyper-parameter. Additionally, as shown in Table 3, when $C = 1$, the performance differs significantly from the base Stable Diffusion. This is because the input logits for the summary token are all replaced with negative infinity.

G.3 HYPER-PARAMETERS IN MEMORIZED IMAGE TRIGGER PROMPT SEARCHING LEVERAGING MCMC

In Table 6, we present the hyper-parameters used in our algorithm for finding memorized image trigger prompts via MCMC.

H ADDITIONAL EXAMPLES OF MEMORIZED IMAGES

In this section, we present the trigger prompts identified by our algorithm along with the generated images from Stable Diffusion using these prompts. Additionally, for each image, we provide the corresponding images presumed to be from the training data, identified using the Reverse Image Search API. The layout repetition of the generated images and those found through the API strongly indicate that Stable Diffusion has memorized the training data. Moreover, we have confirmed that the majority of these images are currently available for commercial sale. We leveraged DDIM (Song et al., 2021a) Scheduler to generate images.

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198



Figure 7: Examples of memorized images found using the Reverse Image Search API. The prompt used for image generation is “Cozy kitchen painted”.

1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212

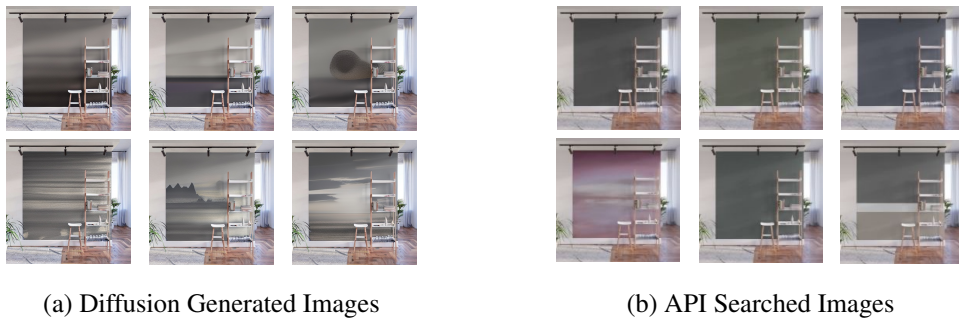


Figure 8: Examples of memorized images found using the Reverse Image Search API. The prompt used for image generation is “Grey standard wall mural”.

1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225



Figure 9: Examples of memorized images found using the Reverse Image Search API. The prompt used for image generation is “Iphone case covered with skull”.

1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239

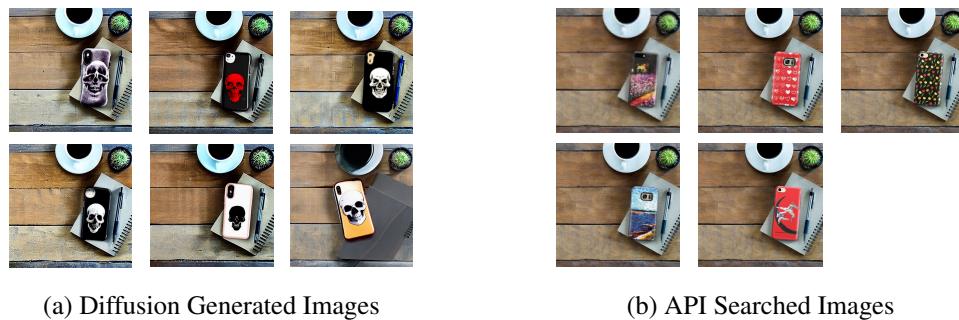


Figure 10: Examples of memorized images found using the Reverse Image Search API. The prompt used for image generation is “Iphone case covered with skull”.

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252



(a) Diffusion Generated Images



(b) API Searched Images

Figure 11: Examples of memorized images found using the Reverse Image Search API. The prompt used for image generation is “Iphone case covered with skull”.

1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266



(a) Diffusion Generated Images



(b) API Searched Images

Figure 12: Examples of memorized images found using the Reverse Image Search API. The prompt used for image generation is “Knit line Africa American quilt house lace boots”.

1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279



(a) Diffusion Generated Images



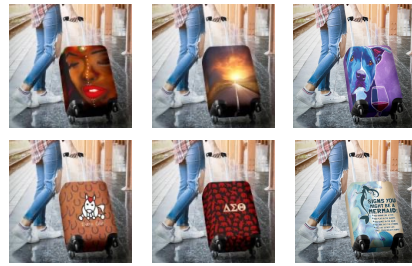
(b) API Searched Images

Figure 13: Examples of memorized images found using the Reverse Image Search API. The prompt used for image generation is “Knit line Africa American quilt house lace boots”.

1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293



(a) Diffusion Generated Images



(b) API Searched Images

Figure 14: Examples of memorized images found using the Reverse Image Search API. The prompt used for image generation is “Travel luggage cover”.

1296

1297

1298

1299

1300

1301

1302

1303

1304

1305

1306

1307

1308



(a) Diffusion Generated Images



(b) API Searched Images

Figure 15: Examples of memorized images found using the Reverse Image Search API. The prompt used for image generation is “United states throw blanket”.

1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

1321

1322

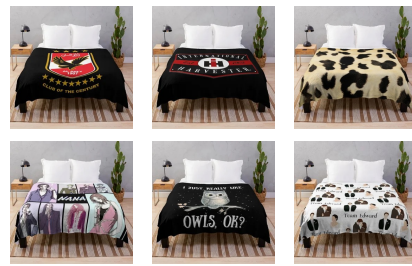
1323

1324

1325



(a) Diffusion Generated Images



(b) API Searched Images

Figure 16: Examples of memorized images found using the Reverse Image Search API. The prompt used for image generation is “United states throw blanket”.

1330

1331

1332

1333

1334

1335

1336

1337

1338

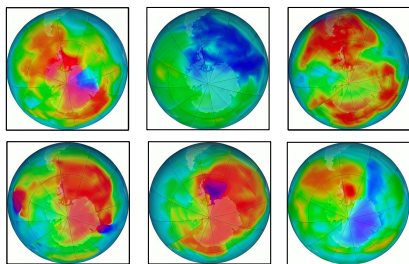
1339

1340

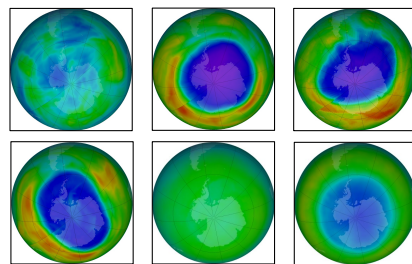
1341

1342

1343



(a) Diffusion Generated Images



(b) API Searched Images

Figure 17: Examples of memorized images found using the Reverse Image Search API. The prompt used for image generation is “Uranus center as ozone temperature map”.

1348

1349

1350 I DATASHEET

1351

1352 I.1 MOTIVATION

1353

1354 **Q1 For what purpose was the dataset created?** Was there a specific task in mind? Was there
1355 a specific gap that needed to be filled? Please provide a description.

1356

1357 • MemBench is a benchmark designed for evaluating memorization mitigation methods
1358 in diffusion models. Recently, many diffusion models have been highlighted for their
1359 issues with image memorization, prompting the development of various memorization
1360 mitigation methods. However, due to the absence of a benchmark to properly evaluate
1361 these methods, their effectiveness has not been adequately assessed. To address this,
1362 we developed MemBench, which includes a large number of memorized image trigger
prompts and appropriate metrics for evaluation.

1363

1364 **Q2 Who created the dataset (e.g., which team, research group) and on behalf of which
1365 entity (e.g., company, institution, organization)?**

1366

1367 • Considering a double-blind review, we will not disclose this information at the current
1368 stage. We will open it to the public in the camera-ready submission.

1369

1370 **Q3 Who funded the creation of the dataset?** If there is an associated grant, please provide the
1371 name of the grantor and the grant name and number.

1372

1373 • Considering a double-blind review, we will not disclose this information at the current
1374 stage. We will open it to the public in the camera-ready submission.

1375

1376 **Q4 Any other comments?**

1377

1378 • No.

1379

1380 I.2 COMPOSITION

1381

1382 **Q5 What do the instances that comprise the dataset represent (e.g., documents, photos,
1383 people, countries)?** *Are there multiple types of instances (e.g., movies, users, and ratings;
1384 people and interactions between them; nodes and edges)? Please provide a description.*

1385

1386 • It includes links to the images memorized by Text-to-Image diffusion models and the
1387 prompts that trigger these images.

1388

1389 **Q6 How many instances are there in total (of each type, if appropriate)?**

1390

1391 • Please refer to Section 1

1392

1393 **Q7 Does the dataset contain all possible instances or is it a sample (not necessarily random)
1394 of instances from a larger set?** *If the dataset is a sample, then what is the larger set? Is the
1395 sample representative of the larger set (e.g., geographic coverage)? If so, please describe
1396 how this representativeness was validated/verified. If it is not representative of the larger set,
1397 please describe why not (e.g., to cover a more diverse range of instances, because instances
1398 were withheld or unavailable).*

1399

1400 • It will be a sample of all existing trigger prompts that induce memorized images in Stable
1401 Diffusion. However, to the best of our knowledge, we have secured the largest number of
1402 trigger prompts, and we plan to add more in the future.

1403

1404 **Q8 What data does each instance consist of?** *“Raw” data (e.g., unprocessed text or images)
1405 or features? In either case, please provide a description.*

1406

1407 • Input trigger prompt and URLs of memorized images triggered by the corresponding
1408 prompt.

1409

1410 **Q9 Is there a label or target associated with each instance?** *If so, please provide a description.*

1411

1412 • No.

1413

1414 **Q10 Is any information missing from individual instances?** *If so, please provide a description,
1415 explaining why this information is missing (e.g., because it was unavailable). This does not
1416 include intentionally removed information, but might include, e.g., redacted text.*

1417

1418 • No.

- 1404 Q11 **Are relationships between individual instances made explicit (e.g., users’ movie ratings,**
1405 **social network links)?** *If so, please describe how these relationships are made explicit.*
1406
- 1407 • Yes, in our benchmark, relationships between individual instances are made explicit. For
1408 example, each trigger prompt in our dataset is explicitly linked to the memorized image
1409 it induces. This is done by including pairs of trigger prompts and their corresponding
1410 memorized images, clearly showing the relationship between them. Additionally, each
1411 image in the benchmark is linked to the specific Text-to-Image diffusion model that
1412 memorized it, providing a clear mapping of model-instance relationships.
- 1413 Q12 **Are there recommended data splits (e.g., training, development/validation, testing)?** *If*
1414 *so, please provide a description of these splits, explaining the rationale behind them.*
1415
- 1416 • No.
- 1417 Q13 **Are there any errors, sources of noise, or redundancies in the dataset?** *If so, please*
1418 *provide a description.*
1419
- 1420 • No.
- 1421 Q14 **Is the dataset self-contained, or does it link to or otherwise rely on external resources**
1422 **(e.g., websites, tweets, other datasets)?** *If it links to or relies on external resources, a) are*
1423 *there guarantees that they will exist, and remain constant, over time; b) are there official*
1424 *archival versions of the complete dataset (i.e., including the external resources as they*
1425 *existed at the time the dataset was created); c) are there any restrictions (e.g., licenses, fees)*
1426 *associated with any of the external resources that might apply to a future user? Please*
1427 *provide descriptions of all external resources and any restrictions associated with them, as*
1428 *well as links or other access points, as appropriate.*
1429
- 1429 • The dataset includes the URLs of the memorized images.
 - 1430 • Regarding (a), we provide multiple URLs for each memorized image to ensure the dataset’s
1431 longevity, even if one hosting source goes down.
 - 1432 • Regarding (b), since the memorized images are not copyrighted by us, we cannot provide
1433 the images directly.
 - 1434 • Regarding (c), these images should be used solely for evaluating the effectiveness of
1435 mitigation methods and should not be used for commercial training or distribution.
- 1436 Q15 **Does the dataset contain data that might be considered confidential (e.g., data that is**
1437 **protected by legal privilege or by doctor–patient confidentiality, data that includes the**
1438 **content of individuals’ non-public communications)?** *If so, please provide a description.*
1439
- 1440 • No.
- 1441 Q16 **Does the dataset contain data that, if viewed directly, might be offensive, insulting,**
1442 **threatening, or might otherwise cause anxiety?** *If so, please describe why.*
1443
- 1444 • No.
- 1445 Q17 **Does the dataset relate to people?** *If not, you may skip the remaining questions in this*
1446 *section.*
1447
- 1448 • No.
- 1449 Q18 **Does the dataset identify any subpopulations (e.g., by age, gender)?**
1450
- 1451 • No.
- 1452 Q19 **Is it possible to identify individuals (i.e., one or more natural persons), either directly or**
1453 **indirectly (i.e., in combination with other data) from the dataset?** *If so, please describe*
1454 *how.*
1455
- 1456 • The memorized images include faces of celebrities, such as Emma Watson.
- 1457 Q20 **Does the dataset contain data that might be considered sensitive in any way (e.g., data**
that reveals racial or ethnic origins, sexual orientations, religious beliefs, political
opinions or union memberships, or locations; financial or health data; biometric or
genetic data; forms of government identification, such as social security numbers;
criminal history)? *If so, please provide a description.*

1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

- No.

Q21 Any other comments?

- Although the memorized images we discovered do not contain offensive or confidential elements, they do include images of currently sold products and faces of celebrities. For instance, there are images of Emma Watson. Therefore, these images should be used solely for evaluating the effectiveness of memorization mitigation methods.

I.3 COLLECTION PROCESS

Q22 How was the data associated with each instance acquired? *Was the data directly observable (e.g., raw text, movie ratings), reported by subjects (e.g., survey responses), or indirectly inferred/derived from other data (e.g., part-of-speech tags, model-based guesses for age or language)? If data was reported by subjects or indirectly inferred/derived from other data, was the data validated/verified? If so, please describe how.*

- Our method discovers memorized images without any prior information by leveraging BERT models, diffusion models, and the Reverse Image Search API. For more details, please refer to Section 3.

Q23 What mechanisms or procedures were used to collect the data (e.g., hardware apparatus or sensor, manual human curation, software program, software API)? How were these mechanisms or procedures validated?

- Our method employs the Reverse Image Search API⁷ and Google Image Search API⁸ to discover memorized images. For more details, please refer to Section 3, 5.

Q24 If the dataset is a sample from a larger set, what was the sampling strategy (e.g., deterministic, probabilistic with specific sampling probabilities)?

- No.

Q25 Who was involved in the data collection process (e.g., students, crowdworkers, contractors) and how were they compensated (e.g., how much were crowdworkers paid)?

- None. The process was automated.

Q26 Over what timeframe was the data collected? Does this timeframe match the creation timeframe of the data associated with the instances (e.g., recent crawl of old news articles)? If not, please describe the timeframe in which the data associated with the instances was created.

- The data was collected from April 2024 to May 2024.

Q27 Were any ethical review processes conducted (e.g., by an institutional review board)? If so, please provide a description of these review processes, including the outcomes, as well as a link or other access point to any supporting documentation.

- No.

Q28 Does the dataset relate to people? If not, you may skip the remaining questions in this section.

- People may appear in the memorized images.

Q29 Did you collect the data from the individuals in question directly, or obtain it via third parties or other sources (e.g., websites)?

- Our method employs the Reverse Image Search API and Google Image Search API to discover memorized images. For more details, please refer to Section 3, 5.

Q30 Were the individuals in question notified about the data collection? If so, please describe (or show with screenshots or other information) how notice was provided, and provide a link or other access point to, or otherwise reproduce, the exact language of the notification itself.

- Our automated memorized image trigger prompt searching algorithm did not involve any participation of individuals.

⁷<https://tineye.com/>

⁸<https://developers.google.com/custom-search>

- 1512 Q31 **Did the individuals in question consent to the collection and use of their data?** *If so,*
1513 *please describe (or show with screenshots or other information) how consent was requested*
1514 *and provided, and provide a link or other access point to, or otherwise reproduce, the exact*
1515 *language to which the individuals consented.*
- 1516 • Our automated memorized image trigger prompt searching algorithm did not involve any
1517 participation of individuals.
- 1518 Q32 **If consent was obtained, were the consenting individuals provided with a mechanism to**
1519 **revoke their consent in the future or for certain uses?** *If so, please provide a description,*
1520 *as well as a link or other access point to the mechanism (if appropriate).*
- 1521 • Our automated memorized image trigger prompt searching algorithm did not involve any
1522 participation of individuals.
- 1523 Q33 **Has an analysis of the potential impact of the dataset and its use on data subjects (e.g.,**
1524 **a data protection impact analysis) been conducted?** *If so, please provide a description*
1525 *of this analysis, including the outcomes, as well as a link or other access point to any*
1526 *supporting documentation.*
- 1527 • We discuss the limitation of our current work in Section 8, and we plan to further investi-
1528 gate and analyze the impact of our benchmark in future work.
- 1529 Q34 **Any other comments?**
- 1530 • No.

1533 I.4 PREPROCESSING, CLEANING, AND/OR LABELING

- 1534 Q35 **Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucket-**
1535 **ing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances,**
1536 **processing of missing values)?** *If so, please provide a description. If not, you may skip the*
1537 *remainder of the questions in this section.*
- 1538 • No.
- 1539 Q36 **Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to**
1540 **support unanticipated future uses)?** *If so, please provide a link or other access point to*
1541 *the “raw” data.*
- 1542 • N/A.
- 1543 Q37 **Is the software used to preprocess/clean/label the instances available?** *If so, please*
1544 *provide a link or other access point.*
- 1545 • N/A.
- 1546 Q38 **Any other comments?**
- 1547 • No.

1551 I.5 USES

- 1552 Q39 **Has the dataset been used for any tasks already?** *If so, please provide a description.*
- 1553 • Not yet. MemBench is a new benchmark.
- 1554 Q40 **Is there a repository that links to any or all papers or systems that use the dataset?** *If*
1555 *so, please provide a link or other access point.*
- 1556 • Not yet. We plan to provide links to works that use our benchmark.
- 1557 Q41 **What (other) tasks could the dataset be used for?**
- 1558 • Image memorization mitigation in diffusion models.
- 1559 Q42 **Is there anything about the composition of the dataset or the way it was collected**
1560 **and preprocessed/cleaned/labeled that might impact future uses?** *For example, is there*
1561 *anything that a future user might need to know to avoid uses that could result in unfair*
1562 *treatment of individuals or groups (e.g., stereotyping, quality of service issues) or other*
1563 *undesirable harms (e.g., financial harms, legal risks) If so, please provide a description. Is*
1564 *there anything a future user could do to mitigate these undesirable harms?*

- 1566
- No.
- 1567
- 1568 Q43 **Are there tasks for which the dataset should not be used?** *If so, please provide a*
- 1569 *description.*
- The images we provide must not be used for training generative models. Since these
 - 1570 images include faces of celebrities and currently sold products, they should never be used
 - 1571 or distributed for the training of generative models.
- 1572
- 1573 Q44 **Any other comments?**
- No.
- 1574
- 1575

1576 I.6 DISTRIBUTION AND LICENSE

- 1577
- 1578 Q45 **Will the dataset be distributed to third parties outside of the entity (e.g., company,**
- 1579 **institution, organization) on behalf of which the dataset was created?** *If so, please*
- 1580 *provide a description.*
- Yes, this benchmark will be open-source.
- 1581
- 1582 Q46 **How will the dataset be distributed (e.g., tarball on website, API, GitHub)?** *Does the*
- 1583 *dataset have a digital object identifier (DOI)?*
- We plan to distribute the formatted data through GitHub after the camera-ready submission.
- 1584
- 1585
- 1586 Q47 **When will the dataset be distributed?**
- After Cam-ready of NeurIPS 2024 dataset and benchmark track.
- 1587
- 1588
- 1589 Q48 **Will the dataset be distributed under a copyright or other intellectual property (IP)**
- 1590 **license, and/or under applicable terms of use (ToU)?** *If so, please describe this license*
- 1591 *and/or ToU, and provide a link or other access point to, or otherwise reproduce, any relevant*
- 1592 *licensing terms or ToU, as well as any fees associated with these restrictions.*
- The dataset will be distributed under the Creative Commons Attribution 4.0 International
 - 1593 (CC BY 4.0) license for the URLs and trigger prompts, not for the images themselves,
 - 1594 as the images themselves are not owned by us. We will provide terms of use document
 - 1595 specifying that the dataset is intended solely for research and evaluation of memorization
 - 1596 mitigation methods and should not be used for training generative models. The GitHub
 - 1597 repository, where the benchmark will be distributed, will contain the code licensed under
 - 1598 the MIT License. The terms of use and licensing information will be accessible via the
 - 1599 GitHub repository when it becomes available.
- 1600
- 1601 Q49 **Have any third parties imposed IP-based or other restrictions on the data associated**
- 1602 **with the instances?** *If so, please describe these restrictions, and provide a link or other*
- 1603 *access point to, or otherwise reproduce, any relevant licensing terms, as well as any fees*
- 1604 *associated with these restrictions.*
- Yes, third parties own the images referenced by the URLs in our dataset. These images
 - 1605 include those of celebrities and currently sold products, which are protected under their
 - 1606 respective intellectual property rights. The URLs provided are for reference purposes only
 - 1607 and must not be used for training or commercial distribution. Any use of the images must
 - 1608 comply with the respective third-party terms and conditions. There are no fees associated
 - 1609 with these restrictions, but users must respect the IP rights of the original content owners.
- 1610
- 1611 Q50 **Do any export controls or other regulatory restrictions apply to the dataset or to**
- 1612 **individual instances?** *If so, please describe these restrictions, and provide a link or other*
- 1613 *access point to, or otherwise reproduce, any supporting documentation.*
- No.
- 1614
- 1615 Q51 **Any other comments?**
- No.
- 1616
- 1617

1618 I.7 MAINTENANCE

- 1619 Q52 **Who will be supporting/hosting/maintaining the dataset?**

- 1620
- 1621
- 1622
- 1623
- 1624
- 1625
- 1626
- 1627
- 1628
- 1629
- 1630
- 1631
- 1632
- 1633
- 1634
- 1635
- 1636
- 1637
- 1638
- 1639
- 1640
- 1641
- 1642
- 1643
- 1644
- 1645
- 1646
- 1647
- 1648
- 1649
- 1650
- 1651
- 1652
- 1653
- 1654
- 1655
- 1656
- 1657
- 1658
- 1659
- 1660
- 1661
- 1662
- 1663
- 1664
- 1665
- 1666
- 1667
- 1668
- 1669
- 1670
- 1671
- 1672
- 1673
- Considering a double-blind review, we will not disclose this information at the current stage. We will open it to the public in the camera-ready submission.
- Q53 How can the owner/curator/manager of the dataset be contacted (e.g., email address)?**
- Through the GitHub discussions that will be opened soon.
 - Through the email of the author.
 - Considering a double-blind review, we will not disclose this information at the current stage. We will open it to the public in the camera-ready submission.
- Q54 Is there an erratum? If so, please provide a link or other access point.**
- No.
- Q55 Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete instances)? If so, please describe how often, by whom, and how updates will be communicated to users (e.g., mailing list, GitHub)?**
- MemBench will be updated. We plan to search for more memorized image trigger prompts and corresponding memorized images using our continuous algorithm.
- Q56 If the dataset relates to people, are there applicable limits on the retention of the data associated with the instances (e.g., were individuals in question told that their data would be retained for a fixed period of time and then deleted)? If so, please describe these limits and explain how they will be enforced.**
- N/A.
- Q57 Will older versions of the dataset continue to be supported/hosted/maintained? If so, please describe how. If not, please describe how its obsolescence will be communicated to users.**
- We will host other versions.
- Q58 If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do so? If so, please provide a description. Will these contributions be validated/verified? If so, please describe how. If not, why not? Is there a process for communicating/distributing these contributions to other users? If so, please provide a description.**
- Through the email of the author.
 - Considering a double-blind review, we will not disclose this information at the current stage. We will open it to the public in the camera-ready submission.
- Q59 Any other comments?**
- No.