EFFICIENT MACHINE UNLEARNING FOR DEEP GENER ATIVE MODELS BY MITIGATING OPTIMIZATION CON FLICTS

Anonymous authors

006

008 009 010

011 012 013

014

015

016

017

018

019

021

025

026

027

028 029

031

032

Paper under double-blind review

Abstract

Machine unlearning of deep generative model refers to the process of modifying or updating a pre-trained generative model to forget or remove certain patterns or information it has learned. Existing research on Bayesian-based unlearning from various deep generative models has highlighted low efficiency as a significant drawback due to two primary causes. Firstly, Bayesian methods often overlook correlations between data to forget and data to remember, leading to conflicts during gradient descent and much slower convergence. Additionally, they require aligning updated model parameters with the original ones to maintain the generation ability of the updated model, further reducing efficiency. To address these limitations, we propose an Efficient Bayesian-based Unlearning method for various deep generative models called EBU. By identifying the relevant weights pertaining to the data to forget and the data to remember, EBU only preserves the parameters related to data to remember, improving the efficiency. Additionally, EBU balances the gradient descent directions of shared parameters to adeptly manage the conflicts caused by the correlations between data to forget and data to remember, leading to a more efficient unlearning process. Extensive experiments on multiple generative models demonstrate the superiority of our proposed EBU.

1 INTRODUCTION

034 In recent years, there have been significant advancements in deep generative models, showcasing their ability to produce synthetic images of exceptional quality (Wei et al., 2022; Li et al., 2022; Nichol 035 & Dhariwal, 2021). These models typically rely on large volumes of training data to effectively 036 learn and generate high-quality outputs (Wang et al., 2022; Cai & Zhu, 2015). However, the use of 037 unauthorized data for training can lead to issues such as data misuse and privacy breaches (Li et al., 2021; Deepa et al., 2022; Wang et al., 2023c), raising concerns about the potential for these models to generate misleading or inappropriate content (Heng & Soh, 2024; Fan et al., 2024). Consequently, 040 there is an urgent need to develop methods for mitigating the influence of specific data on pre-trained 041 generative models. 042

The machine unlearning concept (Bourtoule et al., 2021) is proposed to demonstrate the problem 043 that requires the trained machine learning models *unlearn* from specific data instances. Significant 044 efforts have been made to advance the field of machine unlearning (Gupta et al., 2021; Sekhari et al., 2021; Nguyen et al., 2022). When it comes to the unlearning of multiple deep generative models, prior 046 research has tackled this task by proposing Bayesian-based unlearning methods, albeit with certain 047 inefficiency limitations (Chen et al., 2021; Deepanjali et al., 2021; Schuhmann et al., 2022; Heng & 048 Soh, 2024; Fan et al., 2024). We attribute the inefficiency of Bayesian-based machine unlearning methods to two primary causes. Firstly, they neglect correlations between data to remember and data to forget (Heng & Soh, 2024; Wang et al., 2023a; Fan et al., 2024), leading to conflicts during gradient 051 descent and much slower convergence, thereby exacerbating inherent inefficiencies. Secondly, Bayesian-based unlearning methods require aligning updated model parameters with the original ones 052 to maintain generative ability of updated model (Heng & Soh, 2024; Wang et al., 2023a), introducing additional inefficiencies, notably prolonging the unlearning process.

054 Given the inefficiency limitations of existing Bayesian-based unlearning methods (Chen et al., 2021; 055 Deepanjali et al., 2021; Schuhmann et al., 2022; Fan et al., 2024), we seek to enhance the efficiency of 056 Bayesian-based forgetting methods. We propose EBU to stress the two fold limitations of Bayesian-057 based unlearning methods. Firstly, we selectively retain memory of parameters crucial for data to 058 remember (Sener & Koltun, 2018; Désidéri, 2012), while updating parameters associated with data to forget, preventing the alignment of the entire model parameters, thereby improving efficiency and accelerating the unlearning process. Moreover, considering the correlation between data to forget and 060 data to remember, EBU balances gradient updates on shared parameters associated with both types 061 of data. This balancing mitigates conflicting gradient descent directions, narrowing conflicts, and 062 further enhancing the efficiency of the unlearning process. Notably, the proposed EBU substantially 063 improves the efficiency of the unlearning process and lays the groundwork for more effective model 064 adaptation in deep generative models. 065

066 067

068

069

071

073

075

076

077

078 079

We summary our main contributions in short as follows:

- We introduce EBU, a groundbreaking framework that dramatically enhances the efficiency of machine unlearning in deep generative models by effectively resolving the conflicts between forgetting and remembering processes. This innovation significantly improves both concept-wise and class-wise unlearning.
- EBU strategically preserves critical parameters tied to the data that must be remembered, making the unlearning process in Bayesian-based models far more efficient and streamlined than previous approaches.
- By incorporating a novel mechanism to balance gradient updates for forgetting and remembering, EBU accelerates the entire unlearning process, ensuring faster and more reliable performance.
- Extensive experiments across diverse datasets and generative models clearly demonstrate the superior performance of EBU, proving its effectiveness and efficiency compared to existing baseline methods.
- 081 082 083

RELATED WORK 2

084 085 Machine unlearning for generative model In recent years, the researchers have made efforts in 087 unlearning of generative model. Several works proposed to unlearning from GANs by utilizing the 088 discriminator (Kong & Chaudhuri, 2023; Chen et al., 2021; Sun et al., 2023), but these methods can't be applied to other generative models due to they only suit for paradigm of GANs. There are 089 some works realizing unlearning of generative model by modifying the weights (Bau et al., 2020; Tarun et al., 2023), but it is still a challenge to accurately identify model weights associated to the 091 forgetting tasks accurately. Some researchers proposed Bayesian based unlearning methods that can 092 be applied to various generative models (Heng & Soh, 2024; Nguyen et al., 2020; Fan et al., 2024; Fu et al., 2022; 2021), but they need to trade off between forgetting the posterior distribution of data to 094 forget and not entirely forgetting posterior distribution of the original training data to preserve the 095 generative models' ability, causing conflicts due to the correlations between the data to forget and the 096 original training data.

Difference between EBU and prior methods Our proposed EBU is a Bayesian based unlearning 098 method and can be compatible with any generative models, but advanced with the existing Bayesian based unlearning methods (Heng & Soh, 2024; Nguyen et al., 2020; Fan et al., 2024), our method 100 reduce the conflicts by balancing the gradient descent directions of the parameters shared by the 101 forgetting and remembering processes. We identify these shared parameters by analyzing correspond-102 ing weight saliency maps during the unlearning process. In contrast to previous methods that use 103 saliency maps to identify parameters for updates (Fan et al., 2024), our approach dynamically selects 104 parameters specifically related to forgetting and remembering during the fine-tuning process. This 105 dynamic selection allows us to guide the gradient descent directions for forgetting and remembering separately, while keeping unrelated parameters unchanged. By preventing optimization conflicts 106 between the two tasks, our method not only enhances unlearning efficiency but also ensures more 107 precise parameter updates, leading to superior performance.

¹⁰⁸ 3 PRELIMINARY

110 3.1 PROBLEM DEFINITION

112 We give a brief introduction of the forgetting process of the deep generative model. Given a pre-113 trained generative model G_{θ} with parameter θ trained on the dataset $D = \{(X^n, Y^n)\}_{n=1}^N$, where N is the number of categories in the dataset. Without accessing the training data, we generate forgetting 114 set D_f and remembering set D_r using the model G_{θ} . Here, $D_f = \{(X_f^n, Y_f^n)\}_{n=1}^{N_f}$ denotes the set 115 116 of data to forget and $D_r = \{(X_r^n, Y_r^n)\}_{n=1}^{N_r}$ denotes the set of data to remember, with $N_r + N_f \le N$. In this context, X_f^n and Y_f^n denote the data to forget and corresponding labels for category n. X_r^n 117 118 and Y_r^n denote the data to remember and the corresponding labels for category n. Our goal is to 119 forget the assigned set D_f from the pre-trained generative model G_{θ} while keeping the generation 120 quality of the remaining samples in D_r by fine-turning the pre-trained deep generative model. The 121 fine-tuned model G_{θ^*} with parameters θ^* is expected to forget all samples in D_f while retaining the ability to generate samples conforming to the distribution $p_{\theta^*}(D_r) \sim G_{\theta^*}(X_r|Y_r)$ that is expected 122 to align with the distribution of the data to remember $p(D_r)$. For clarity and ease of representation, 123 all model parameters used in this paper are denoted as the set of their individual elements. For 124 instance, $\theta = \{\theta_1, \theta_2, \dots, \theta_i, \dots, \theta_k\}$, where k represents the total number of elements within θ and 125 θ_i denotes the i^{th} element of the model parameters. 126

127 128

3.2 MOTIVATION

129 We are interested in forgetting specified samples from a pre-trained diffusion model. Prior work 130 attempted to forget data from a model to ensure the privacy in machine learning models by deleting 131 only specific shards (Bourtoule et al., 2021), thereby forgetting these assigned shards. Moreover, 132 building on the ideas from (Heng & Soh, 2024; Nguyen et al., 2020), they implemented forgetting 133 of model by forgetting the posterior belief of data D_f while not forgetting the posterior belief 134 given the full data D. As suggested by (Heng & Soh, 2024), using Elastic Weight Consolidation 135 (EWC) (Kirkpatrick et al., 2017) keeps the posterior belief of the full data $p_{\theta}(D)$, preventing 136 catastrophic forgetting. However, this approach hinders the forgetting process by maintaining the posterior belief of full data $p_{\theta}(D)$. A more reasonable solution is to remember the posterior belief of 137 D_r while forgetting the posterior belief of D_f . The weights for forgetting $p_{\theta}(D_f)$ and remembering 138 $p_{\theta}(D_r)$ are denoted as θ_f and θ_r , respectively. Fast forgetting of D_f can be achieved by keeping θ_r 139 consistent with its original values and leaving θ_f unchanged. Moreover, the overarching unlearning 140 process necessitates a delicate trade-off between forgetting $p_{\theta}(D_f)$ while retaining aspects of $p_{\theta}(D)$. 141 The parameters affected by this trade-off are those within $\theta_f \cap \theta = \theta_f$. When compared to the 142 general unlearning approach, focusing solely on addressing conflicts within $\theta_r \cap \theta_f \subseteq \theta_f$ results in 143 a smaller negative impact on the unlearning process, thereby potentially accelerating it.

144 145

4 Method

146 147 148

149

150

151

152

153

We introduce a novel unlearning method EBU aimed at expediting the forgetting process of pretrained deep generative models while preserving the quality of the generated images. We present a comprehensive theoretical analysis to underpin our proposed EBU in Section 4.1, elucidating the underlying rationale behind our method's efficacy in resolving conflicts inherent in the unlearning process. In particular, EBU comprises two key components. First, we introduce a partial parameter alignment method in Section 4.2, which significantly enhances unlearning efficiency. Second, in Section 4.3, we propose an effective approach to mitigate optimization conflicts between remembering and forgetting, further accelerating the unlearning process.

154 155 156

157

4.1 THEORETICAL ANALYSIS

The optimization objective of the unlearning of deep generative models is to minimize the expected loss functions \mathcal{L}_f and \mathcal{L}_r over the distributions of forgetting dataset $p(D_f)$ and remembering dataset $p(D_r)$:

$$\mathcal{O}_1 = \min_{\{\boldsymbol{\theta}_f, \boldsymbol{\theta}_r\}} \mathbb{E}_{p(D_f)} [\mathcal{L}_f(X_f, \boldsymbol{\theta}_f)] + \mathbb{E}_{p(D_r)} [\mathcal{L}_r(X_r, \boldsymbol{\theta}_r)]$$
(1)

But the distributions $p(D_f)$ and $p(D_r)$ are unavailable, and one solution is to optimize the forgetting task and remembering task independently. However, it's important to note that the data to forget and the data to remember might be related. If handle the remembering and forgetting separately, we could lose important information because the parameters involved in both tasks might conflict with each other (Sener & Koltun, 2018). Thus training the data of forgetting and remembering simultaneously is another solution, which results in a new optimization task:

$$\mathcal{O}_2 = \min_{\boldsymbol{\theta}} \frac{1}{N_f + N_r} \sum_{i=1}^{N_f} \mathcal{L}_f(X_f^i, \boldsymbol{\theta}) + \sum_{i=1}^{N_r} \mathcal{L}_r(X_r^i, \boldsymbol{\theta})$$
(2)

However, this training approach overlooks the distinction between the data intended for forgetting 171 and the data intended for remembering. As a result, it may hinder the model from converging since 172 the gradient of the first loss term and the second loss term in Equation (2) may update in opposite 173 direction. In the context of deep generative model, according to the PAC-Bayesian theory (McAllester, 174 1998), there exists a error $\epsilon(N,\zeta) \ge 0$ with probability $1-\zeta$ over independent draws monotonically 175 decreasing with the training samples N between the expected optimization objective and actual 176 optimization objective (McAllester, 1998) (the detailed proof is provided in Appendix C). We can 177 obtain the following proposition: 178

Proposition 4.1. The bounds between the optimization objective \mathcal{O}_1 and \mathcal{O}_2 :

$$|\mathcal{O}_1 - \mathcal{O}_2| \leq \mathbb{E}_{p(D_f)}[G_{\theta_f}(X_f|Y_f) - G_{\theta_e}(X_f|Y_f)] + \mathbb{E}_{p(D_r)}[G_{\theta_r}(X_r|Y_r) - G_{\theta_e}(X_r|Y_r)] + \epsilon(N_f + N_r, \zeta)$$
(3)

where θ_e denotes the optimal solution of Equation (2) and $G_{\theta_{(\cdot)}}$ denotes the deep generative model with parameters $\theta_{(\cdot)}$.

Considering the differences and connections between the data to forget and the data to remember at the same time, we rewrite the optimization objective as:

187 188

189 190

191

198

199

202

203

204 205

185

179 180 181

168 169 170

$$\mathcal{O}^* = \min_{\{\boldsymbol{\theta}_f, \boldsymbol{\theta}_r, \boldsymbol{\theta}_e\}} \frac{1}{N_f} \sum_{i=1}^{N_f} \mathcal{L}_f(X_f^i, \boldsymbol{\theta}_f) + \frac{1}{N_r} \sum_{i=1}^{N_r} \mathcal{L}_f(X_r^i, \boldsymbol{\theta}_r)$$
subject to $|\mathbb{E}_{p(D_f)}[G_{\boldsymbol{\theta}_f}(X_f|Y_f) - G_{\boldsymbol{\theta}_e}(X_f|Y_f)]| \le \xi$
(4)

$$|\mathbb{E}_{p(D_r)}[G_{\boldsymbol{\theta}_r}(X_r|Y_r) - G_{\boldsymbol{\theta}_e}(X_r|Y_r)]| \le \xi$$

The constant ξ controls the closeness between functions. A larger ξ allows functions to be more task-specific. The expectations in Equation (4) can't be calculated directly due to the lack of the accessibility to the probability distribution $p(\cdot)$. But if the function is Lipschitz in the parameterization, the distance between the functions can be measured by the distance between parameters (Cervino et al., 2021), thus we have another proposition to estimate the expectations:

Proposition 4.2. The two generative models can be seen as two parametric functions, thus it has:

$$\mathbb{E}_{p(D_f)}[G_{\theta_f}(X_f|Y_f) - G_{\theta_e}(X_f|Y_f)] \le L|\theta_f - \theta_e|, \\
\mathbb{E}_{p(D_r)}[G_{\theta_r}(X_r|Y_r) - G_{\theta_e}(X_r|Y_r)] \le L|\theta_r - \theta_e|$$
(5)

where *L* is a constant that decides the scope.

Thus through imposing constraints on parameters, the Equation (4) can be solved. The forgetting and remembering tasks relate to part parameters θ_f and θ_r , we present how to select and modify the corresponding θ_f and θ_r to unlearn from deep generative model in Section 4.2.

206 4.2 PARTIAL PARAMETER ALIGNED BAYESIAN UNLEARNING

207 Existing works (Chen et al., 2021; Deepanjali et al., 2021; Schuhmann et al., 2022) for unlearning 208 of deep generative models need to *fully align* the parameters between the unlearned model and 209 original pre-trained model. This leads to highly inefficient unlearning. In this section, we propose 210 a technique which only needs to *partially align* the parameters between the unlearned model and 211 original pre-trained model, thereby improving the efficiency significantly. Specifically, in Section 4.1, 212 we define the optimization objective of unlearning from deep generative models in Equation (4), but 213 the expectations can't be calculated directly. As described in Proposition 4.2, the expectations can be estimated by accessing the corresponding parameters, then this optimization objective can be solved. 214 In this section, we propose efficient ways to select the corresponding parameters θ_f and θ_r , thereby 215 solving the Equation (4) to unlearn from the deep generative models efficiently.

235 236

245 246

252

253 254

260

Learning objective for the data to forget To forget D_f from the pre-trained deep generative model G_{θ} , we need to make the posterior distribution $p(X_f|\theta^*, Y_f)$ far from the real distribution $p(X_f|Y_f)$ as much as possible. One way is to make the unlearned model to generate samples from a mandatory distribution $\tilde{q}(X_f|Y_f)$, which should be different from the real distribution (Heng & Soh, 2024) and usually set it to a standard Gaussian distribution. We minimize the KL divergence between the generated sample distribution from the unlearned model and the designated mandatory distribution during the fine-turning process:

$$\mathcal{L}_f = \mathbb{E}_{p(Y_f)} D_{KL}(p(X_f | \boldsymbol{\theta}^*, Y_f) || \tilde{q}(X_f | Y_f))$$
(6)

we denote the class distributions in D_f as $p(Y_f)$. In this way, the fine-tuned generative model G_{θ^*} will change the learned distribution of the forgotten set to *mandatory distribution* after fine-tuning, so as to achieve the purpose of making the model forget. Accordingly, the parameters θ_f related to forgetting task should be far from the original parameters.

Learning objective for the data to remember During the forgetting process, it's necessary to strengthen the memory of generative model for the data in the remember set D_r . We thus replay the data from D_r during the forgetting process to keep the model's ability to generate the data to remember. By replaying data in D_r , the updated posterior distribution of D_r denoted as $p(X_r | \theta^*, Y_r)$ is forced to approach to the original posterior distribution $p(X_r | \theta, Y_r)$, preserving the model's ability to generate the data to remember. To replay data from D_r , we define an optimization function as:

$$\mathcal{L}_r = \mathbb{E}_{p(Y_r)} D_{KL}(p(X_r | \boldsymbol{\theta}^*, Y_r) || p(X_r | \boldsymbol{\theta}, Y_r))$$
(7)

through this optimization function, the generated sample distribution of D_r with the unlearned model G_{θ^*} will be consistent with the distribution of the data to remember on the original model G_{θ} .

During the optimization process, the posterior distribution of D_r on model G_{θ^*} should be close to the posterior distribution on G_{θ} , whereas on D_f the situation is reverse. That is we should keep the model parameters θ_r related to remembering close to the corresponding parameters in θ and θ_f related to forgetting far from the corresponding parameters in θ . We use θ_r^{-f} to represent the remaining parameters of θ_r after removing the overlap with the θ_f , and $\hat{\theta}_r^{-f}$ denotes the elements of the original pre-trained model parameters θ with same parameter index of the elements in θ_r^{-f} :

$$\hat{\boldsymbol{\theta}}_r^{-f} = \boldsymbol{\theta} \times \mathcal{O}(\boldsymbol{\theta}_r^{-f}), \ \boldsymbol{\theta}_r^{-f} = \boldsymbol{\theta}_r \setminus (\boldsymbol{\theta}_r \cap \boldsymbol{\theta}_f)$$
(8)

where $\mathcal{O}(\cdot)$ sets the non-zero values in the set to 1 and leave the zero values unchanged, the symbol × denotes element-wise multiplication and \ denotes the operation of obtaining complement set.

Motivated by the continue learning algorithm EWC (Kirkpatrick et al., 2017), we make θ_r close to θ at each gradient step to enforce the model remembering D_r , and the optimization function of remembering D_r can be further rewritten as follows,

J

$$\mathcal{L}_r^c = \mathcal{L}_r + \frac{\gamma}{2} \sum_i F_i (\boldsymbol{\theta}_{r,i}^{-f} - \hat{\boldsymbol{\theta}}_{r,i}^{-f})^2 \tag{9}$$

where γ is a constant that regularizes the degree that make the new parameters close to the old parameters, *i* denotes the index of the element, and *F* is the set of diagonal elements of the fisher information matrix calculated on θ . Existing work (Chen et al., 2021; Heng & Soh, 2024) aligns the entire parameters θ and θ^* . In contrast, our work only needs to partially align $\theta_{r,i}^{-f}$ and $\hat{\theta}_{r,i}^{-f}$ in Equation (9), thereby speeding up the unlearning process.

261 4.3 Optimization Conflicts Mitigation

262 Forgetting data D_f from the generative model changes part of the parameters, while replaying data 263 D_r to the model also changes part of parameters, we denote the two parts of parameters as θ_f and 264 θ_r respectively. However, data from different categories may exhibit similar patterns, potentially 265 influencing a shared set of model parameters, *i.e.*, $\theta_f \cap \theta_r$. As shown in Figure 1, the changed 266 parameters θ_f and θ_r partially overlap. When optimizing \mathcal{L}_r and \mathcal{L}_f simultaneously, the gradient update directions for forgetting and replaying conflict, hindering the model's ability to forget or 267 remember effectively and slowing convergence. To address this, a trade-off strategy is needed 268 to balance forgetting and remembering during optimization. In this section, we propose efficient 269 methods to mitigate these conflicts and accelerate the unlearning process.



(a) Existing unlearning methods (Nguyen et al., (b) Our proposed unlearning method EBU 2020; Heng & Soh, 2024)

Figure 1: In Figure 1(a), we illustrate the conflict inherent in existing unlearning methods, where the conflict spans across all θ_f when forgetting specific data while not entirely forgetting the entire dataset. Figure 1(b) depicts our proposed unlearning method, which narrows the scope of conflicted parameters to $\theta_f \cap \theta_r$, we also employ a trade-off parameter α to balance the conflicted gradient descent directions, mitigating their negative impact on the unlearning process.

Parameters selection During the unlearning process, the parameters that are closely related to the forgetting and remembering tasks will exhibit larger gradients compared to the irrelevant parameters. This observation allows us to identify the parameters pertinent to each task based on their gradients at each gradient step. We denote the entire set of parameters of the deep generative model as θ_t^* at the *t*-th gradient step. The optimization functions for forgetting and remembering at this step are represented by \mathcal{L}_f^t and \mathcal{L}_r^t , respectively. By analyzing these gradients, we can effectively discern which parameters need to be adjusted for efficient unlearning.

We utilize a constant $\sigma \in [0, 1]$ to determine the proportion of selected parameters based on the gradient values $\nabla_{\theta_t^*} \mathcal{L}_f^t$. By applying this threshold, we can identify the parameters θ_f^t associated with forgetting at the *t*-th step by selecting the top σ proportion of gradient values.

$$\boldsymbol{\theta}_{f}^{t} = \boldsymbol{\theta}_{t}^{*} \times \mathcal{O}(\operatorname{top-}\sigma(\nabla_{\boldsymbol{\theta}_{t}^{*}}\mathcal{L}_{f}^{t}))$$
(10)

We fill zeros in the empty spaces to maintain the size of the selected parameters consistent with the original parameters. This ensures that the overall structure of the parameter set remains intact while allowing us to focus on the relevant parameters associated with the forgetting task. By preserving the dimensions of the parameter set, we can seamlessly integrate these updates into the model without disrupting its architecture.

Gradient Modulation To effectively resolve conflicts arising from shared parameters, we propose a
 straightforward trade-off method that balances the effects of forgetting and remembering by averaging
 their gradients. This approach is particularly effective because it neutralizes the competing influences
 of both processes, preventing one from undermining the other.

By defining $\theta_f^{t,r} = \theta_f^t \cap \theta_r^t$ to represent the parameters in θ_f^t that overlap with those in θ_r^t , we ensure that we focus on the shared parameters that are crucial for both forgetting and remembering. Averaging the gradients for these overlapping parameters allows us to update them in a way that accommodates the requirements of both tasks simultaneously:

316

282

284

285

286

287

288

300

306

$$\nabla_{\boldsymbol{\theta}_{f}^{t,r}} \mathcal{L}_{f}^{t} \coloneqq \alpha \nabla_{\boldsymbol{\theta}_{f}^{t,r}} \mathcal{L}_{r}^{t} \oplus (1-\alpha) \nabla_{\boldsymbol{\theta}_{r}^{t,f}} \mathcal{L}_{f}^{t}$$
(11)

where \oplus denotes the element-wise sum operation applied to elements with the same index, while where \oplus denotes the element-wise sum operation applied to elements with the same index, while are $\alpha \in (0, 1)$ serves as a trade-off constant. Adjusting α allows us to steer the gradient descent directions of the parameters shared between θ_f^t and θ_r^t (we delve into the effect of different α values on the unlearning process in Section 5.5).

This method effectively mitigates optimization conflicts by ensuring that updates made for forgetting do not completely override the updates for remembering, and vice versa. As a result, we maintain a balanced influence on the shared parameters, which leads to a more coherent unlearning process. This dual consideration enhances overall model efficiency and ensures that important information is preserved while unwanted knowledge is successfully removed.

During the optimization process, it is crucial to effectively erase the data we wish to forget from the deep generative model. To achieve this, we adjust the gradient descent to be more focused on forgetting, thereby accelerating the forgetting process. We denote $\mathcal{L}_r^{c,t}$ as the *t*-th instance of the forgetting loss \mathcal{L}_r^c (refer to Equation Equation (9)).

 $\nabla_{\boldsymbol{\theta}_{t}^{t}}(\mathcal{L}_{f}^{t} + \mathcal{L}_{r}^{c,t}) := \nabla_{\boldsymbol{\theta}_{t}^{t}}(\mathcal{L}_{f}^{t} + \mathcal{L}_{r}^{c,t}) + \nabla_{\boldsymbol{\theta}_{f}^{t}}\mathcal{L}_{f}^{t}$ (12)

This targeted approach ensures that the model prioritizes updates that facilitate the removal of unwanted information while maintaining the integrity of the data we intend to remember.

Consequently, the overall optimization process can be summarized as:

$$\boldsymbol{\theta}_{t+1}^* \leftarrow \boldsymbol{\theta}_t^* - \lambda \nabla_{\boldsymbol{\theta}_t^*} (\mathcal{L}_f^t + \mathcal{L}_r^{c,t}), \quad \boldsymbol{\theta}_0^* = \boldsymbol{\theta}$$
(13)

where λ is the learning rate. This formulation highlights how the parameters are updated by taking into account both the forgetting and remembering objectives, ensuring a balanced approach that facilitates effective unlearning. The comprehensive algorithm of our proposed method is outlined in Algorithm 1.

Algorithm 1 Unlearning Process of EBU

Input: Pre-trained model G_θ, data to retain D_r = {(Xⁿ_r, Yⁿ_r)}^{N_r}_{n=1}, data to forget D_f = {(Xⁿ_f, Yⁿ_f)}^{N_f}_{n=1}, desired distribution q̃(X_f|Y_f), learning rate λ, initial parameter set θ^{*}₀ = θ
 Output: Fine-tuned model G_{θ*}, θ* = θ^{*}_T
 for t = 0 to T-1 do
 Compute gradients Σ_C C^t = ∂C^t / ∂θ* and Σ_C C^t = ∂C^t / ∂θ*

- 4: Compute gradients $\nabla_{\boldsymbol{\theta}_{t}^{*}} \mathcal{L}_{f}^{t} = \partial \mathcal{L}_{f}^{t} / \partial \boldsymbol{\theta}_{t}^{*}$ and $\nabla_{\boldsymbol{\theta}_{t}^{*}} \mathcal{L}_{r}^{t} = \partial \mathcal{L}_{r}^{t} / \partial \boldsymbol{\theta}_{t}^{*}$
- 5: Select parameter subsets θ_f^t and θ_r^t based on the top δ proportion values of $\nabla_{\theta_t^*} \mathcal{L}_f^t$ and $\nabla_{\theta_t^*} \mathcal{L}_r^t$ respectively.
- 6: Identify overlapping parameters: $\boldsymbol{\theta}_{f}^{t,r} = \boldsymbol{\theta}_{r}^{t,f} = \boldsymbol{\theta}_{r}^{t} \cap \boldsymbol{\theta}_{r}^{t}$
- 7: Calculate gradients for overlapping parameters: $\nabla_{\boldsymbol{\theta}_{f}^{t,r}} \mathcal{L}_{f}^{t} := \alpha \nabla_{\boldsymbol{\theta}_{f}^{t,r}} \mathcal{L}_{f}^{t} \oplus (1-\alpha) \nabla_{\boldsymbol{\theta}_{r}^{t,f}} \mathcal{L}_{r}^{t}$
- 8: Update parameters: $\boldsymbol{\theta}_{t+1}^* \leftarrow \boldsymbol{\theta}_t^* \lambda \nabla_{\boldsymbol{\theta}_t^*} (\mathcal{L}_f^t + \mathcal{L}_r^c), \quad \nabla_{\boldsymbol{\theta}_t^*} (\mathcal{L}_f^t + \mathcal{L}_r^c) \leftarrow \nabla_{\boldsymbol{\theta}_t^*} (\mathcal{L}_f^t + \mathcal{L}_r^c) + \nabla_{\boldsymbol{\theta}_t^*} \mathcal{L}_f^t$
- 9: end for
- 355 356 357

358

360

361

362

331 332

333

334

335 336 337

338

339

340

341 342

343

344

345

347

348

349

350

351

352

353

354

5 EXPERIMENT

In this section, we demonstrate the ability of proposed EBU in assisting various deep generate models unlearning certain classes and concepts. We compare our method with the existing state-of-the-art unlearning baselines, highlighting the effectiveness of EBU.

5.1 EXPERIMENT SETTING

Implements We focus on two types of forgetting tasks: class-wise forgetting and concept-wise forgetting. We utilize two types of generative models: the pre-trained DDPM (Ho et al., 2020) and the Stable Diffusion (SD) model (Rombach et al., 2022) to assess the performance of our proposed method. The hyperparameters α and δ are set to 0.6 and 0.5, respectively. All experiments are conducted using 4 Nvidia V100 GPUs with 32 GB memory. More detailed experiment implements of both the class-wise unlearning and concept-wise unlearning tasks can be found in Appendix B.2.

Baselines We compare our proposed EBU with other five different state-of-the-art methods to evaluate the efficiency and fidelity of EBU. We choose three general unlearning methods: FT (Warnecke et al., 2021), GA (Thudi et al., 2022) and Retraining, two unlearning methods for deep generative methods: SA (Heng & Soh, 2024) and ESD (Gandikota et al., 2023). The detailed description and implements of the baseline methods are presented in Appendix B.

376 Metrics To evaluate the fidelity of unlearning methods for class-wise forgetting, we use two metrics:
 377 classification entropy (CE) for forgetting classes and remaining accuracy (RA) for the accuracy of the remaining classes in the unlearned model. For assessing efficiency, we consider unlearn time

(UT) and relearn time (RT). UT measures the time for the unlearning process, while RT counts the gradient updating steps needed for the unlearned model. We also use the Fréchet Inception Distance (FID) to assess the image quality of the classes to remember, aiming for minimal impact on their quality by the unlearned model. As for the concept-wise forgetting, we use two metrics: Clip Score (CS) (Hessel et al., 2021) and Nudity Score (NS) used in (Gandikota et al., 2023) to evaluate the forgetting performance of different unlearning methods.

5.2 **CLASS-WISE FORGETTING**

The class-wise forgetting is to unlearn the specified classes, we evaluate the class-wise forgetting performance of unlearning methods on pre-trained DDPM and SD.

Main results on DDPM We conduct experiments on CIFAR10, STL10 and CIIFAR100. The experi-ment results of DDPM on CIFAR-100 and STL10 are presented in Table 6, the experiment results on CIFAR10 are presented in Appendix D.4, multiple unlearning methods are applied to pre-trained DDPM to demonstrate the effectiveness of unlearning methods. Here we present the experiment results of unlearning class 0, and the additional experiment results are shown in Appendix D.

Table 1: The experiment results of different unlearning methods on CIFAR100 and STL10 datasets with pre-trained DDPM, and class 0 is selected to be unlearned. The best results are bolded and the second best results are underlined.

| Method | CIFAR-100 | | | | | STL10 | | | | |
|----------------------------|------------------|-------------------------------------|------------------------------------|--------------------------------------|--------------------|-----------------------------------|-------------------------------------|------------------------------------|------------------------------------|-----------------------------|
| | $CE_f(\uparrow)$ | $\operatorname{CE}_r({\downarrow})$ | $\mathrm{RA}\left(\uparrow\right)$ | $\mathrm{RT}\left(\downarrow\right)$ | $FID~(\downarrow)$ | $\operatorname{CE}_{f}(\uparrow)$ | $\operatorname{CE}_r({\downarrow})$ | $\mathrm{RA}\left(\uparrow\right)$ | $\text{RT}\left(\downarrow\right)$ | $\text{FID} \ (\downarrow)$ |
| FT (Warnecke et al., 2021) | 1.247 | 1.946 | 0.245 | 2000 | 298.6 | 1.631 | 1.628 | 0.104 | 2000 | 187.4 |
| GA (Thudi et al., 2022) | 1.222 | 1.499 | 0.239 | 2000 | 40.44 | 1.749 | 1.722 | 0.120 | 2000 | 332.7 |
| SA (Heng & Soh, 2024) | 1.306 | 1.463 | 0.401 | 20000 | 40.28 | 1.822 | 0.089 | 0.968 | 30000 | 48.87 |
| SalUn (Fan et al., 2024) | 1.218 | 1.482 | 0.381 | 2000 | 59.28 | 0.596 | 0.092 | 0.990 | 4000 | 75.91 |
| EBU (Ours) | 1.431 | 1.398 | 0.419 | 200 | 37.88 | 1.917 | 0.086 | 0.955 | 200 | 48.35 |

It can be seen from Table 1 that our EBU demonstrates a significant reduction in RT while maintaining performance on par with other baseline methods. Notably, our approach achieves superior results across most evaluation metrics, emphasizing both its efficiency and robustness. Moreover, EBU excels in CIFAR100 datasets with a higher number of classes, where the intricate correlations between the data to forget and the data to remember are more effectively managed. This highlights our method's ability to handle the optimization conflicts between the remembering and forgetting, improving the efficiency and ensuring reliable and scalable unlearning performance.



(a) CE of D_f varies with unlearning time.



Figure 2: The changes of cross-entropy (CE) of D_f and D_r with different unlearning methods when unlearning time increases.

To further demonstrate the efficiency of our proposed method, we plotted the changes in classification cross-entropy of CIFAR10 when unlearning class 0 during the unlearning process in Figure 5 (more details are presented in Appendix D.2). It's evident that our method can forget the designated categories within a short time (within 10^2 seconds) without affecting the remaining data. Moreover, to further validate the presence of optimization conflicts on shared parameters, we report the average
 number of parameters related to both forgetting and remembering during the fine-tuning process in
 Appendix D.1.

Main results on SD We perform the class-wise forgetting on pre-trained standard SD model, ten classes of Imagenette are chosen to evaluate the performance of unlearning from SD. We present the experiment results of forgetting class 'cassette player' in Table 2, and the detailed experiment settings can be found in Appendix B. We further present the generated samples of different unlearning methods with prompt "An image of cassette player" in Figure 9 (refer to Appendix D.4).

It can be seen from Figure 9 that, 441 compared with SA, our proposed 442 method drastically reduces the time re-443 quired for forgetting, with only 1000 444 steps required to achieve good re-445 sults whereas SA requires 50000 steps 446 for complete forgetting, resulting in 447 $50 \times$ efficiency improvement, also 448 compared with ESD and SalUn, our 449 method has better unlearning perfor-

Table 2: Experiment results of different unlearning methods on pre-trained SD, class 'cassette player' is selected to be unlearned. The best results are bolded and the second best results are underlined.

| Mathad | Imagenette | | | | | | | | |
|-------------------------------|-----------------------------------|-------------------------------------|----------------------------|------------------------------------|-----------------------------|--|--|--|--|
| Method | $\operatorname{CE}_{f}(\uparrow)$ | $\operatorname{CE}_{r}(\downarrow)$ | RA (†) | $\text{RT}\left(\downarrow\right)$ | FID (\downarrow) | | | | |
| ESD (Gandikota et al., 2023). | $1.089_{\pm.120}$ | $0.159_{\pm.022}$ | $0.936_{\pm.010}$ | 1000 | $201.9_{+2.33}$ | | | | |
| SA (Heng & Soh, 2024) | $1.196_{\pm.084}$ | $0.027_{\pm.002}$ | $0.998_{\pm,001}$ | 50000 | $211.7_{\pm 4.01}$ | | | | |
| SalUn (Fan et al., 2024) | $1.139_{\pm.024}$ | $0.054_{\pm .032}$ | $0.976_{\pm.011}$ | 1000 | $201.7_{\pm 2.11}$ | | | | |
| EBU (Ours) | $ 1.148_{\pm.054} $ | $\underline{0.041}_{\pm.005}$ | $\textbf{0.999}_{\pm.001}$ | 1000 | $\textbf{199.9}_{\pm 2.01}$ | | | | |

mance with same gradient descent steps and better preserves the model's ability to generate samples
 of data to remember. The experiment results demonstrate the effectiveness and efficiency of our
 proposed methods in unlearning from deep generative models.

Unlearning process visualization We visualize the unlearning process of forgetting class 0 from
 the pre-trained DDPM in Figure 6 (refer to Appendix D.3), displaying a total of 200 steps. The
 samples of all baseline methods can be found in Appendix D. Our proposed method exhibits superior
 forgetting performance with the fewest forgetting time steps.

457

459

458 5.3 CONCEPT-WISE FORGETTING

Concept-wise forgetting involves the unlearning of specific concepts, often employed in the text-to-image models. In this study, we evaluate the concept-wise forgetting performance of various unlearning methods on the classic text-to-image model, SD. Our methodology begins by generating samples with empty prompts. Subsequently, we establish the *mandatory distribution* of samples for specific concepts as the distribution of these randomly generated samples.

Forgetting Nudity We assess the effectiveness
of our EBU in forgetting nudity, the quantitative results are presented in Table 3. Moreover,
we illustrate the performance of unlearning nudity in Figure 3 (more samples can refer to Appendix D.5). To ensure fair comparison across

| Table 3: The experiment results of unlearning "nu- | - |
|----------------------------------------------------|---|
| dity" on SD model, four baseline methods are used | |

| Metric | EBU | SA | ESD | SalUn | SPM |
|--------------|--------|--------|--------|--------|--------|
| Clip Score | 0.1900 | 0.1747 | 0.1895 | 0.1396 | 0.1874 |
| Nudity Score | 0.5988 | 0.4615 | 0.3557 | 0.5062 | 0.5528 |

experiments, we set the gradient descent steps to 1000 for all unlearning methods. Using the prompt "a person with full nudity," we generate samples with different seeds.





471

Figure 3: The generated samples of nude person with prompt "a person with full nudity".

Comparative analysis with baseline methods shows that our approach consistently achieves the
 lowest levels of nudity, even though the number of gradient descent steps remains the same. This
 demonstrates the effectiveness of our method in selectively unlearning undesired content while
 maintaining efficiency. We also evaluate the performance of forgetting art style in Appendix D.5.

486 5.4 ABLATION STUDY

488 We have also conducted an ablation 489 study of our proposed EBU. To eval-490 uate the effectiveness of the Partial 491 Align (PA) module, we remove the PA module (denoted as $EBU_{w/o PA}$). 492 As for the effectiveness of Optimiza-493 tion Conflicts Mitigation module, we 494 remove the \mathcal{L}_f , \mathcal{L}_r^c and the mitiga-495 tion operation respectively. The abla-496 tion study is performed on CIFAR10 497 dataset and the results are presented 498 in Table 4. As shown in Table 4, the 499

| Table 4: | The | ablation | study | of | our | EBU | method | on | the |
|----------|-------|------------|-------|----|-----|-----|--------|----|-----|
| CIFAR-1 | 0 dat | aset using | g DDP | M. | | | | | |

| Ablation | $CE_f(\uparrow)$ | $\operatorname{CE}_{r}\left(\downarrow\right)$ | RA (†) | $\mathrm{RT}\left(\downarrow\right)$ | FID (\downarrow) | | | | | |
|-----------------------------------|------------------|------------------------------------------------|--------|--------------------------------------|--------------------|--|--|--|--|--|
| Partial Align | | | | | | | | | | |
| $\mathrm{EBU}_{w/o \mathrm{PA}}$ | 0.8517 | 0.0353 | 0.9625 | 200 | 33.78 | | | | | |
| Optimization Conflicts Mitigation | | | | | | | | | | |
| $EBU_{w/o \mathcal{L}_f}$ | 0.8213 | 0.0224 | 0.9888 | 200 | 35.38 | | | | | |
| $EBU_{w/o \mathcal{L}_{x}^{c}}$ | 0.8212 | 0.0227 | 0.9886 | 200 | 35.37 | | | | | |
| EBUw/o Mitigation | 0.8033 | 0.0386 | 0.9888 | 200 | 42.42 | | | | | |
| EBU | 0.8550 | 0.0110 | 0.9931 | 200 | 29.92 | | | | | |

removal of the Proximal Attention (PA) significantly degrades the forgetting performance, underscoring its critical role. Furthermore, the contributions of both the forgetting loss \mathcal{L}_f and the reconstruction loss \mathcal{L}_r are evident in improving the overall forgetting. The proposed optimization conflict mitigation mechanism effectively reduces conflicts between objectives, leading to enhanced forgetting performance. These results demonstrate that both PA and conflict mitigation are essential components for optimizing the forgetting process.

5.5 Effect of α and δ

We assess the impact of α and δ on the forgetting process by varying their values within the range 0.2, 0.4, 0.6, 0.8. Our experiments are conducted on CIFAR10, and results are depicted in Figure 4. **The effect of** α : We observe that increasing α affects the generated samples of both data to remember



Figure 4: We investigate the impact of α and δ on the forgetting performance through experiments conducted on the CIFAR10 dataset. We vary the values of α and δ within the range 0.2, 0.4, 0.6, 0.8.

and data to forget simultaneously. Specifically, within a small range of α , the quality of generated samples for data to remember improves. Conversely, within a larger range, the performance of forgetting is enhanced. **The effect of** δ : Furthermore, the value of δ determines the proportion of parameters attributed to forgetting and remembering. As δ increases, the quality of generated samples for data to remember improves. However, when δ reaches a large value, the forgetting process is impacted adversely due to the involvement of more irrelevant parameters.

533

500

501

502

503

504 505 506

507 508

509

510

511

521

522 523 524

526

527

528

6 CONCLUSION AND LIMITATION

In conclusion, we addressed the inefficiencies existed in Bayesian-based unlearning methods for deep
 generative models. We proposed an Efficient Bayesian-based Unlearning method (EBU), which significantly enhances the unlearning process. By pinpointing relevant parameters and balancing gradient
 descent directions of data to forget and data to remember, EBU preserves essential parameters and
 manages conflicts effectively, resulting in a more efficient unlearning process. Extensive experiments
 across various generative models and unlearning tasks demonstrate the superior performance of EBU,
 validating its effectiveness and efficiency in unlearning tasks.

540 REFERENCES 541

547

551

552

553

554

567

568

569

570

580

581

582

583

- David Bau, Steven Liu, Tongzhou Wang, Jun-Yan Zhu, and Antonio Torralba. Rewriting a deep 542 generative model. In Computer Vision-ECCV 2020: 16th European Conference, Glasgow, UK, 543 August 23–28, 2020, Proceedings, Part I 16, pp. 351–369. Springer, 2020. 544
- Thomas Baumhauer, Pascal Schöttle, and Matthias Zeppelzauer. Machine unlearning: Linear filtration 546 for logit-based classifiers. Machine Learning, 111(9):3203-3226, 2022.
- Lucas Bourtoule, Varun Chandrasekaran, Christopher A Choquette-Choo, Hengrui Jia, Adelin Travers, 548 Baiwu Zhang, David Lie, and Nicolas Papernot. Machine unlearning. In 2021 IEEE Symposium 549 on Security and Privacy (SP), pp. 141-159. IEEE, 2021. 550
 - Li Cai and Yangyong Zhu. The challenges of data quality and data quality assessment in the big data era. Data science journal, 14:2–2, 2015.
- Yinzhi Cao and Junfeng Yang. Towards making systems forget with machine unlearning. In 2015 IEEE symposium on security and privacy, pp. 463–480. IEEE, 2015. 555
- Olivier Catoni. A pac-bayesian approach to adaptive classification. preprint, 840:2, 2003. 556
- Juan Cervino, Juan Andrés Bazerque, Miguel Calvo-Fullana, and Alejandro Ribeiro. Multi-task 558 reinforcement learning in reproducing kernel hilbert spaces via cross-learning. IEEE Transactions 559 on Signal Processing, 69:5947–5962, 2021. 560
- Kongyang Chen, Yao Huang, and Yiwen Wang. Machine unlearning via gan. arXiv preprint 561 arXiv:2111.11869, 2021. 562
- 563 Natarajan Deepa, Quoc-Viet Pham, Dinh C Nguyen, Sweta Bhattacharya, B Prabadevi, Thippa Reddy 564 Gadekallu, Praveen Kumar Reddy Maddikunta, Fang Fang, and Pubudu N Pathirana. A survey 565 on blockchain for big data: Approaches, opportunities, and future directions. Future Generation 566 Computer Systems, 131:209–226, 2022.
 - S Deepanjali, S Dhivya, and S Monica Catherine. Efficient machine unlearning using general adversarial network. In Artificial Intelligence Techniques for Advanced Computing Applications: Proceedings of ICACT 2020, pp. 487–494. Springer, 2021.
- Jean-Antoine Désidéri. Multiple-gradient descent algorithm (mgda) for multiobjective optimization. 571 *Comptes Rendus Mathematique*, 350(5-6):313–318, 2012. 572
- 573 Chongyu Fan, Jiancheng Liu, Yihua Zhang, Dennis Wei, Eric Wong, and Sijia Liu. Salun: Em-574 powering machine unlearning via gradient-based weight saliency in both image classification and 575 generation. In International Conference on Learning Representations, 2024. 576
- Daniel L Felps, Amelia D Schwickerath, Joyce D Williams, Trung N Vuong, Alan Briggs, Matthew 577 Hunt, Evan Sakmar, David D Saranchak, and Tyler Shumaker. Class clown: Data redaction in 578 machine unlearning at enterprise scale. arXiv preprint arXiv:2012.04699, 2020. 579
 - Shaopeng Fu, Fengxiang He, Yue Xu, and Dacheng Tao. Bayesian inference forgetting. arXiv preprint arXiv:2101.06417, 2021.
 - Shaopeng Fu, Fengxiang He, and Dacheng Tao. Knowledge removal in sampling-based bayesian inference. arXiv preprint arXiv:2203.12964, 2022.
- 585 Rohit Gandikota, Joanna Materzynska, Jaden Fiotto-Kaufman, and David Bau. Erasing concepts 586 from diffusion models. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 2426–2436, 2023.
- 588 Aditya Golatkar, Alessandro Achille, and Stefano Soatto. Eternal sunshine of the spotless net: 589 Selective forgetting in deep networks. In Proceedings of the IEEE/CVF Conference on Computer 590 Vision and Pattern Recognition, pp. 9304–9312, 2020.
- Varun Gupta, Christopher Jung, Seth Neel, Aaron Roth, Saeed Sharifi-Malvajerdi, and Chris Waites. 592 Adaptive machine unlearning. Advances in Neural Information Processing Systems, 34:16319– 16330, 2021.

| 594 595 596 | Alvin Heng and Harold Soh. Selective amnesia: A continual learning approach to forgetting in deep generative models. <i>Advances in Neural Information Processing Systems</i> , 36, 2024. |
|--------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| 597 598 | Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference- free evaluation metric for image captioning. <i>arXiv preprint arXiv:2104.08718</i> , 2021. |
| 599 600 601 | Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. Advances in neural information processing systems, 33:6840–6851, 2020. |
| 602 603 604 605 | James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. <i>Proceedings of the national academy of sciences</i> , 114 (13):3521–3526, 2017. |
| 606 607 608 | Zhifeng Kong and Kamalika Chaudhuri. Data redaction from pre-trained gans. In 2023 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML), pp. 638–677. IEEE, 2023. |
| 609 610 611 | Junyi Li, Tianyi Tang, Wayne Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. Pretrained language models for text generation: A survey. <i>arXiv preprint arXiv:2201.05273</i> , 2022. |
| 612 613 614 | Qinbin Li, Zeyi Wen, Zhaomin Wu, Sixu Hu, Naibo Wang, Yuan Li, Xu Liu, and Bingsheng He. A survey on federated learning systems: Vision, hype and reality for data privacy and protection. <i>IEEE Transactions on Knowledge and Data Engineering</i> , 2021. |
| 615 616 617 | David A McAllester. Some pac-bayesian theorems. In <i>Proceedings of the eleventh annual conference</i> on Computational learning theory, pp. 230–234, 1998. |
| 618 619 | Quoc Phong Nguyen, Bryan Kian Hsiang Low, and Patrick Jaillet. Variational bayesian unlearning. Advances in Neural Information Processing Systems, 33:16025–16036, 2020. |
| 620 621 622 623 | Thanh Tam Nguyen, Thanh Trung Huynh, Phi Le Nguyen, Alan Wee-Chung Liew, Hongzhi Yin, and Quoc Viet Hung Nguyen. A survey of machine unlearning. <i>arXiv preprint arXiv:2209.02299</i> , 2022. |
| 624 625 | Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In <i>International conference on machine learning</i> , pp. 8162–8171. PMLR, 2021. |
| 626 627 628 629 | Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High- resolution image synthesis with latent diffusion models. In <i>Proceedings of the IEEE/CVF confer-</i> <i>ence on computer vision and pattern recognition</i> , pp. 10684–10695, 2022. |
| 630 631 632 633 | Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. <i>Advances in Neural Information Processing Systems</i> , 35:25278–25294, 2022. |
| 634 635 636 637 | Ayush Sekhari, Jayadev Acharya, Gautam Kamath, and Ananda Theertha Suresh. Remember what you want to forget: Algorithms for machine unlearning. <i>Advances in Neural Information Processing Systems</i> , 34:18075–18086, 2021. |
| 638 639 | Ozan Sener and Vladlen Koltun. Multi-task learning as multi-objective optimization. Advances in neural information processing systems, 31, 2018. |
| 640 641 642 | Hui Sun, Tianqing Zhu, Wenhan Chang, and Wanlei Zhou. Generative adversarial networks unlearn- ing. <i>arXiv preprint arXiv:2308.09881</i> , 2023. |
| 643 644 645 | Ayush K Tarun, Vikram S Chundawat, Murari Mandal, and Mohan Kankanhalli. Fast yet effective machine unlearning. <i>IEEE Transactions on Neural Networks and Learning Systems</i> , 2023. |
| 646 647 | Anvith Thudi, Gabriel Deza, Varun Chandrasekaran, and Nicolas Papernot. Unrolling sgd: Under- standing factors influencing machine unlearning. In 2022 IEEE 7th European Symposium on Security and Privacy (EuroS&P), pp. 303–319. IEEE, 2022. |

| 0.40 | |
|-------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| 648 649 650 | Shuohang Wang, Yichong Xu, Yuwei Fang, Yang Liu, Siqi Sun, Ruochen Xu, Chenguang Zhu, and Michael Zeng. Training data is more valuable than you think: A simple and effective method by retrieving from training data. <i>arXiv preprint arXiv:2203.08773</i> , 2022. |
| 651 | |
| 652 | Weiqi Wang, Zhiyi Tian, Chenhan Zhang, An Liu, and Shui Yu. Bfu: Bayesian federated unlearning |
| 653 | with parameter self-sharing. In Proceedings of the 2023 ACM Asia Conference on Computer and |
| 654 | Communications Security, pp. 567–578, 2023a. |
| GEE | |
| 055 | Weiqi Wang, Chenhan Zhang, Zhiyi Tian, and Shui Yu. Machine unlearning via representation |
| 656 | forgetting with parameter self-sharing. <i>IEEE Transactions on Information Forensics and Security</i> , |
| 657 | 2023b. |
| 658 | Then Wang, Farming Vang, Li Shan, and Hang, Huang, A comprehensive survey of forgetting in |
| 659 | does howing bound optimus housing a Win proprint arVin 2207 00218, 2002 |
| 660 | deep learning beyond continual learning. arxiv preprint arxiv:2507.09218, 2025c. |
| 661 | Zhenyi Wang, Yan Li, Li Shen, and Heng Huang. A unified and general framework for continual |
| 662 | learning. In The Twelfth International Conference on Learning Representations, 2024. |
| 663 | |
| 664 | Alexander Warnecke, Lukas Pirch, Christian Wressnegger, and Konrad Rieck. Machine unlearning |
| 665 | of features and labels. arXiv preprint arXiv:2108.11577, 2021. |
| 666 | Jacon Wai Vi Tay Dichi Dommesoni Colin Doffel Derret Zonh Schootien Derregoud Dani Vegetame |
| 667 | Jason wei, 11 Tay, Kisin Dominasani, Comin Kaner, Barlet Zopii, Sebastian Bolgeauu, Dani Togatania, |
| 660 | arViv preprint arViv:2206.07682, 2022 |
| 000 | <i>urxiv preprint urxiv.</i> 2200.07082, 2022. |
| 669 | Peng-Fei Zhang, Guangdong Bai, Zi Huang, and Xin-Shun Xu. Machine unlearning for image |
| 670 | retrieval: A generative scrubbing approach. In <i>Proceedings of the 30th ACM International</i> |
| 671 | Conference on Multimedia, pp. 237–245, 2022. |
| 672 | |
| 673 | |
| 674 | |
| 675 | |
| 676 | |
| 677 | |
| 678 | |
| 679 | |
| 680 | |
| 681 | |
| 682 | |
| 683 | |
| 684 | |
| 685 | |
| 686 | |
| 607 | |
| 007 | |
| 688 | |
| 689 | |
| 690 | |
| 691 | |
| 692 | |
| 693 | |
| 694 | |
| 695 | |
| 696 | |
| 697 | |
| 698 | |
| 699 | |
| 700 | |
| 701 | |

A RELATED WORK OF MACHINE UNLEARNING

704 Machine unlearning dedicates to enabling the removal of specific data from trained machine learning 705 models (Bourtoule et al., 2021; Nguyen et al., 2022), encompasses methods ranging from data 706 reorganization to direct model manipulation. While some approaches focus on preventing model 707 learning by reorganizing data or constructing fake data (Felps et al., 2020; Tarun et al., 2023; Zhang 708 et al., 2022; Cao & Yang, 2015), they often introduce extra information that impacts learning. Some researchers tried to unlearn sample through manipulating the model directly (Baumhauer et al., 2022; 709 Golatkar et al., 2020; Sekhari et al., 2021; Wang et al., 2024; 2023b), i.e., modifying weights to remove 710 sample information or adjusting model updates based on data statistics during training (Golatkar 711 et al., 2020; Sekhari et al., 2021). Notably, research primarily focuses on unlearning for classification 712 models, with relatively less attention to generative models, which is still in its early stages. 713

714 715

B EXPERIMENT SETTING

- 716
- 717 718

719 B.1 THE DESCRIPTION OF BASELINE METHODS

General unlearning methods: We choose Retraining, FT (Warnecke et al., 2021) and GA (Thudi et al., 2022) as basic unlearning baseline methods. Retraining realizes unlearning by updating on the training data with removal of data. FT builds on close-form updates of model parameters to unlearn the features and labels. GA limits the overall change in weights during SGD to facilitate the approximate unlearning.

Unlearning for generative models: We use two most recent unlearning methods for generative models SA (Heng & Soh, 2024) and ESD (Gandikota et al., 2023). SA derives from continual learning to selectively forget concepts in pretrained deep generative models. ESD erases a visual concept from a pre-trained diffusion model, given only the name of the style and using negative guidance as a teacher model.

731 732

B.2 IMPLEMENT DETAILS

Class-wise forgetting: Class-wise forgetting targets the removal of generations belonging to specified classes from deep generative models. For the DDPM, we fine-tune it for 200 gradient update steps using a batch size of 32 and a learning rate of 1*e*-5. We conduct comparisons on two datasets, CIFAR-10 and STL10, employing various methods. Regarding the SD model, we fine-tune it for 1000 gradient update steps with a learning rate of 1*e*-5. We focus on unlearning classes from Imagenette, which comprises ten easily identifiable classes from ImageNet. The detailed experiment results can be found in Section 5.2.

Concept-wise forgetting: The concept-wise unlearning of deep generative models aims to eliminate generations containing specific concepts. We fine-tune the SD model using 1000 gradient descent steps, with a batch size set to 1 and a learning rate of 1*e*-5. We consider two types of concepts: art style and nudity. Detailed experiments and results can be found in Section 5.3.

744 745

746

B.3 THE DETAILS OF FORGETTING FROM DDPM

Baseline implement. For the baseline methods, we implement them as recommend. Note that the baseline methods FT, GA are proposed for classification model, but they can be adapted to generative model easily. For FT, it only needs data to remember during the unlearning process, thus we change the original loss function of FT as \mathcal{L}_r . For GA, it only needs data to forget, thus we set the loss function of GA as $-\mathcal{L}_f$. And for Retraining, we just fine-tune the pre-trained model with D_r directly. For SA, we implement it as recommend.

Experiment settings In our experimental setup, we utilize a simple yet effective Residual model architecture with 3 input and output channels, employing a channel size of 128 and 2 residual blocks.
Attention resolutions are set at 16, with dropout probability at 0.1. For diffusion, we implement a linear beta schedule spanning from 0.0001 to 0.02 over 1000 diffusion timesteps. During training, we

employ a batch size of 32 for 20,0 iterations, with logging every 50 iterations and visualization of 100 samples. Optimization is conducted using the Adam optimizer with a learning rate of 0.0001 and a weight decay of 0.000, while gradients are clipped at a threshold of 1.0. These settings ensure robust experimentation and reliable evaluation of our proposed methods.

B.4 THE EXPERIMENT SETTINGS OF FORGETTING FROM SD.

Baseline implement We use two baseline methods ESD and SD. For SD, it needs to generate the data to forget and data to remember, we use the random samples generated by SD as the data to forget and we then use a empty prompt to generate the data to remember. For ESD, we implement it as recommended.

Figure 1767
 Experiment settings In our experiment setup, we utilize the Latent Diffusion model with specific configurations tailored for unlearning tasks. The diffusion process spans 1000 timesteps, with linear beta scheduling from 0.00085 to 0.012. We employ a UNet model architecture with attention resolutions at [4, 2, 1] and two residual blocks. Training involves a base learning rate of 1.0e-05, and we utilize a LambdaLinear scheduler with a warm-up period of 1 step. The model consists of a first stage autoencoder with embedded dimensions of 4 and a conditional stage encoder. And all the unlearning methods are trained with 'xattn' part parameters while keeping other parts frozen.

C THEORY PROOF

C.1 PAC ASSUMPTION

To give the bound between the solution of multi-task learning problem, the following assumptions need to be introduced in advance:

Assumption C.1. If function G_{θ} is probably approximately correct, for all $\theta \in \Theta$ with probability $1 - \delta$ over independent draws $(X_n, Y_n) \sim p$:

$$|\mathbb{E}[\mathcal{L}(X,\boldsymbol{\theta})] - \frac{1}{N} \sum_{i=1}^{N} \mathcal{L}(X_n,\boldsymbol{\theta})| \le \epsilon(N,\delta)$$
(14)

Assumption C.2. Loss function $\mathcal{L}(., \theta)$ is M-Lipschitz continuous.

The Assumption C.1 is a generalization of the law of large numbers for the case in which samples are id, where the error is of under 1/N. And under the assumption, we could obtain the proposition:

Proposition C.3. The bounds between the optimization objective \mathcal{O}_1 and \mathcal{O}_2 :

$$|\mathcal{O}_1 - \mathcal{O}_2| \leq \mathbb{E}_{p(D_f)}[G_{\theta_f}(X_f|Y_f) - G_{\theta_e}(X_f|Y_f)] + \mathbb{E}_{p(D_r)}[G_{\theta_r}(X_r|Y_r) - G_{\theta_e}(X_r|Y_r)] + \epsilon(N_f + N_r, \zeta)$$

$$(15)$$

where θ_e denotes the optimal solution of Equation (2) and $G_{\theta_{(.)}}$ denotes the deep generative model with parameters $\theta_{(.)}$.

The Equation (15) in the proposition is a direct application of Assumption C.1 over the average
 probability distribution and taking the Lipschitz Assumption C.2 over the solutions.

Note that if the function is Lipschitz in the parameterization, there is a connection between the functional, and parametric constraints:

Proposition C.4. The two generative models can be seen as two parametric functions, thus it has:

$$\mathbb{E}_{p(D_f)}[G_{\boldsymbol{\theta}_f}(X_f|Y_f) - G_{\boldsymbol{\theta}_e}(X_f|Y_f)] \leq L|\boldsymbol{\theta}_f - \boldsymbol{\theta}_e|, \\
\mathbb{E}_{p(D_r)}[G_{\boldsymbol{\theta}_r}(X_r|Y_r) - G_{\boldsymbol{\theta}_e}(X_r|Y_r)] \leq L|\boldsymbol{\theta}_r - \boldsymbol{\theta}_e|$$
(16)

Through enforcing the constraint over the parameters, we could remove the expectation and the dependency over the distribution $p(\cdot)$.

C.2 THE PAC-BAYESIAN THEORY

PAC-Bayesian theory (McAllester, 1998) seeks to quantify the trade-off between empirical risk minimization and model complexity, offering insights into the generalization ability of a learning algorithm. Given a hypothesis class \mathcal{H} , the training set $D = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$ iid sampled from an distribution \hat{p} over an instance space \mathcal{S} , the real-valued loss function $\mathcal{L}: \mathcal{H} \times \mathcal{S} \longrightarrow [0, \infty)$, the PAC-Bayes provides the generalization bounds for any posterior $q \in \mathcal{M}^1_+(\mathcal{H})$, and $\mathcal{M}^1_+(\mathcal{H})$ is the set of probability measures on a space \mathcal{H} . The generalization bounds are dependent on the empirical performance of q and its closeness to a chosen prior distribution p, the empirical risks of a posterior distribution q are defined as:

$$\mathcal{R}_s(q) = \mathbb{E}_{h \sim q(h)}\left[\frac{1}{n} \sum_{i=1}^n \mathcal{L}(h, (X_i, Y_i))\right]$$
(17)

and the true risk of q is defined as:

$$\mathcal{R}(q) = \mathbb{E}_{h \sim q(h)}[\mathbb{E}_{(X,Y) \sim \hat{p}}\mathcal{L}(h, (X,Y))]$$
(18)

Also before introducing the PAC-Bayesian bound, the posterior q and the loss function \mathcal{L} need to satisfy the Assumption C.5.

Assumption C.5. If there exists a constant K>0 and a family \mathcal{E} of functions $\mathcal{H} \to \mathbb{R}$, for any $(X_1, Y_1), (X_2, Y_2) \in \mathcal{S}:$

$$d_{\mathcal{E}}(q(h|(X_1, Y_1)), q(h|(X_2, Y_2))) \ge Kd((X_1, Y_1), (X_2, Y_2))$$
(19)
and the loss function $\mathcal{L}(., (X, Y)) : \mathcal{H} \to \mathbb{R}$ is in \mathcal{E} .

Then for the bounded loss functions, the PAC-Bayesian bound (Catoni, 2003) are defined in Theo-rem C.6

Theorem C.6. For a probability measure \hat{p} on S, a loss function $\mathcal{L} : \mathcal{H} \times S \longrightarrow [0,1]$, with probability at least 1- δ over the n random samples \hat{S} draw from \hat{p} , the following equation holds for any posterior distribution $q \in \mathcal{M}^1_+(H)$:

$$\mathcal{R}(q) \le \mathcal{R}_s(q) + \frac{\lambda}{8n} + \frac{D_{KL}(q||p) + \log\frac{1}{\delta}}{\lambda}$$
(20)

the real number $\delta \in (0, 1)$ and $\lambda > 0$.

The Theorem C.6 predicts the behavior of q(h|(X,Y)) for any $(X,Y) \sim \hat{p}$ when the posterior q was learned using the training samples \mathcal{D} .

D **EXPERIMENT RESULTS**

D.1 THE COUNT OF OVERLAPPED PARAMETERS

We have counted the average number of selected forgetting and remembering parameters of the SD model, as well as the overlapping parameters during the unlearning process. In our experiments,

Table 5: The average count of parameters related to forgetting and remembering and the overlapped parameters of SD model during the unlearning process.

| Forgetting concept | Total | Forgetting | remembering | overlapped | |
|----------------------------|--------|------------|-------------|------------|--|
| Multiple Artists | 3.200G | 1.605G | 1.988G | 1.095G | |
| Single Artist | 3.200G | 1.701G | 1.596G | 0.867G | |
| Multiple violence concepts | 3.200G | 1.616G | 2.010G | 1.100G | |
| Single violence concepts | 3.200G | 1.675G | 1.566G | 0.861G | |

we observed that there is indeed a non-negligible overlap between the parameters θ_r (retain) and θ_f (forget), and we present the count of overlapped parameters in Table 5. This overlap was identified through the gradient analysis during the training process, where conflicting updates to shared parameters were detected.

D.2 EFFICIENCY ANALYSIS

To further demonstrate the efficiency of our proposed method, we plotted the changes in classification cross-entropy of CIFAR10 when unlearning class 0 during the unlearning process in Figure 5. It's evident from the plot that our method can forget the designated categories within a short time (within 10^2 seconds) without affecting the remaining data. In contrast, other methods consistently require more than 10^3 seconds for the same task. This evidence highlights that our method achieves a significant improvement in unlearning time and is more efficient than other baseline methods.



Figure 5: The changes of cross-entropy (CE) of D_f and D_r with different unlearning methods when unlearning time increases, experiments are performed on CIFAR10 dataset, class 0 is unlearned.

D.3 UNLEARNING PROCESS

We present the visualization of unlearning process of the unlearning methods with 200 gradient descent steps on CIFAR10 dataset in figure 6.

D.4 CLASS-WISE FORGETTING

Forgetting from pre-trained DDPM. The full experiment results on CIFAR10 dataset are presented in Table 6:

Table 6: The experiment results of different unlearning methods on CIFAR10 and STL10 datasets with pre-trained DDPM, and class 0 is selected to be unlearned. The best results are bolded and the second best results are underlined.

| Mathad | | CIFAR-10 | | | | STL10 | | | | | |
|----------------------------|--------------------|-------------------------------------|--------------------|------------------------------------|--------------------------------|-----------------------------------|-----------------------------------|--------------------|------------------------------------|--------------------------------|--|
| Method | $CE_f(\uparrow)$ | $\operatorname{CE}_{r}(\downarrow)$ | RA (†) | $\text{RT}\left(\downarrow\right)$ | FID (\downarrow) | $\operatorname{CE}_{f}(\uparrow)$ | $\operatorname{CE}_r(\downarrow)$ | RA (†) | $\text{RT}\left(\downarrow\right)$ | FID (\downarrow) | |
| FT (Warnecke et al., 2021) | $0.579 \pm .006$ | $0.580_{\pm.001}$ | $0.114_{\pm.001}$ | 2000 | $75.51_{\pm 21.1}$ | $ 1.631_{\pm.162} $ | $1.628_{\pm.191}$ | $0.104_{\pm.005}$ | 2000 | $187.4_{\pm 21.2}$ | |
| GA (Thudi et al., 2022) | $0.627_{\pm.005}$ | $0.609 \pm .005$ | $0.596 \pm .008$ | 2000 | 253.6 ± 15.0 | $1.749 \pm .009$ | $1.722_{\pm.007}$ | $0.120_{\pm.024}$ | 2000 | $332.7_{\pm 10.1}$ | |
| SA (Heng & Soh, 2024) | $0.807_{\pm.006}$ | $0.009_{\pm.001}$ | $0.997_{\pm,001}$ | 20000 | $19.11_{\pm 2.41}$ | $1.822_{+.284}$ | $0.089_{\pm.005}$ | $0.968_{\pm,013}$ | 30000 | $48.87_{\pm 2.81}$ | |
| SalUn (Fan et al., 2024) | $0.598_{\pm.012}$ | $0.061 \pm .004$ | $0.996_{\pm.003}$ | 4000 | $29.91_{\pm 2.12}$ | $0.596 \pm .018$ | $0.092_{\pm .006}$ | $0.990_{\pm.015}$ | 4000 | $75.91_{\pm 2.83}$ | |
| EBU (Ours) | $0.855_{\pm .051}$ | $0.011_{\pm .001}$ | $0.993_{\pm .001}$ | 200 | $\underline{29.92}_{\pm 4.21}$ | $ 1.917_{\pm.021} $ | $\underline{0.086}_{\pm.015}$ | $0.955_{\pm .021}$ | 200 | $\underline{48.35}_{\pm 4.32}$ | |

The visualizations of different unlearning methods unlearn from pre-trained DDPM on CIFAR10 and STL10 are presented on figure 7 and figure 8.

Forgetting from pre-trained SD. We present the generated samples of different unlearning methods with prompt "An image of cassette player" in Figure 9.

- D.5 CONCEPT-WISE FORGETTING
- Forgetting of Nudity. We illustrate the performance of unlearning nudity of four baseline methods in Figure 10.



Figure 6: Samples generated by our method, SA, and FT for the classes 0 ("airplane") and class 1 ("car") on the CIFAR10 dataset. The class to forget is "airplane", and the class to remember is "car". The unlearning step 't' varies from 10 to 200.

Forgetting Art style To assess the efficacy of various forgetting methods in unlearning art styles, we conduct experiments using the SD model. In Figure 11, we present the forgetting results for "Kelly Mckernan" and "Thomas Kinkade". Our EBU demonstrates the capability to effectively forget art styles from the generative model. The figures generated by our method exhibit different colors and objects compared to the original figures generated by the SD model. This showcases the effectiveness of EBU in forgetting art styles.

 D.6 EFFECT OF HYPER-PARAMETER

The impact of α and δ on the forgetting process, their values are chosen in the range 0.2, 0.4, 0.6, 0.8. Our experiments are conducted on CIFAR10, and results are depicted in Figure 12.







Figure 12: We investigate the impact of α and δ on the forgetting performance through experiments conducted on the CIFAR10 dataset. We vary the values of α and δ within the range 0.2, 0.4, 0.6, 0.8.