

# ADVERSARIAL MACHINE UNLEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

This paper focuses on the challenge of machine unlearning, aiming to remove the influence of specific training data on machine learning models. Traditionally, the development of unlearning algorithms runs parallel with that of membership inference attacks (MIA), a type of privacy threat to determine whether a data instance was used for training. However, the two strands are intimately connected: one can view machine unlearning through the lens of MIA success with respect to removed data. Recognizing this connection, we propose a game-theoretic framework that integrates MIAs into the design of unlearning algorithms. Specifically, we model the unlearning problem as a Stackelberg game in which an unlearner strives to unlearn specific training data from a model, while an auditor employs MIAs to detect the traces of the ostensibly removed data. Adopting this adversarial perspective allows the utilization of new attack advancements, facilitating the design of unlearning algorithms. Our framework stands out in two ways. First, it takes an adversarial approach and proactively incorporates the attacks into the design of unlearning algorithms. Secondly, it uses implicit differentiation to obtain the gradients that limit the attacker’s success, thus benefiting the process of unlearning. We present empirical results to demonstrate the effectiveness of the proposed approach for machine unlearning.

## 1 INTRODUCTION

The enactment of the General Data Protection Regulation (GDPR) by the EU has elevated the importance of deleting user data from machine learning models to a critical level. This process is distinctly more intricate compared to removing data from conventional databases. Erasing the data’s imprint from a machine learning model necessitates an approach to negate the data’s influence on the model comprehensively while maintaining the utility and accuracy of the model.

Beyond this, establishing the true extent to which data influence has been erased from the model poses a significant challenge (Song & Mittal, 2021). Numerous methods and metrics have been advanced to validate the thoroughness of data removal, each with varying degrees of reliability and efficacy (Guo et al., 2020; Thudi et al., 2022b). We propose a novel adversarial perspective on unlearning that we argue is a more robust framework for effective machine unlearning. In this approach, the focus shifts to simulating possible attacks aimed at inferring whether the data that should have been forgotten nevertheless maintains some influence on the model. If, within this adversarial framework, an attacker fails to distinguish whether a data point was part of the training set or merely a typical instance of unseen data, we can conclude that the influence of the data point on the data has been successfully unlearned.

We leverage advancements from the burgeoning domain of Membership Inference Attacks (MIA) to simulate an adversary (Shokri et al., 2017), therein framing a Stackelberg Game (SG) between an unlearner, tasked with orchestrating the unlearning process, and an auditor deploying MIA to deduce the membership of data in the model’s training set. The key idea is for the unlearner to adjust the model being unlearned by utilizing gradient feedbacks from the auditor’s optimization problem, moving the model in a direction that limits the effectiveness of the attack, thus achieving the goal of unlearning. Specifically, we formulate the MIA as a utility-maximizing problem, where the utility measures the remaining influence of a data point in the unlearned model. The unlearner’s loss function is defined as a combination of the degradation of model performance and the auditor’s utility. We harness the development from implicit differentiation and design a gradient-based algorithm to

solve the game, allowing for seamless integration into existing end-to-end pipelines (Gould et al., 2016; Amos & Kolter, 2017; Agrawal et al., 2019).

The contributions of the present paper are summarized below

- We propose to evaluate the effectiveness of an unlearning algorithm from an adversarial perspective, inspiring us to develop a game theory framework that enables the use of advanced MIAs for enhancing the unlearning process.
- Additionally, we design a gradient-based solution method to solve the game by leveraging implicit differentiation, making it amenable to end-to-end pipelines.
- Finally, we support the efficacy of the game and the solution method with extensive results.

## 2 RELATED WORK

The first related thread is machine unlearning, which focuses on removing the influence of a subset of data (referred to as the forget set) from a machine learning model. The unlearning approaches are divided into two classes. The first one is exact unlearning, which involves retraining the model on data excluding the forget set. The second one is approximate unlearning. The ideas behind approximate unlearning are twofold. The first is to track the influence of each training data on the updates to a model’s weights, allowing for a rollback during unlearning (Bourtoule et al., 2021; Graves et al., 2021b; Chen et al., 2022). The second is using a loss function to capture the objectives of unlearning (e.g., removing the influence of the forget set while maintaining model utility) and modifying the model weights to minimize the loss function (Guo et al., 2020; Golatkar et al., 2020b; Izzo et al., 2021b; Warnecke et al., 2023; Chundawat et al., 2023; Jia et al., 2023). The method proposed in this paper aligns with the second idea. Specifically, we design a loss function that simulates an auditor who uses MIAs to evaluate the effectiveness of unlearning. By differentiating through the auditor’s optimization problem, we compute the gradients that reduce the auditor’s utility, thus increasing the effectiveness of unlearning. Besides algorithmic developments, Jagielski et al. (2023) proposes a measure to quantify the forgetting during training; Thudi et al. (2022b) takes a formal analysis of the definition of approximation unlearning and propose methods to verify exact unlearning. Due to space constraint, it is not feasible to provide a comprehensive review of all related studies. We refer the readers to the survey article by Nguyen et al. (2022) for a more exhaustive discussion.

The second related line is membership inference attacks (MIA). Shokri et al. (2017) introduced MIAs, showing the privacy risks of machine learning models. Subsequently, different attack methods are proposed (Chen et al., 2021; Carlini et al., 2022; Ye et al., 2022; Bertran et al., 2023). On the other hand, Carlini et al. (2022) shows that existing criteria to evaluate MIAs are limited in capturing real-world scenarios and propose more practical evaluation metrics. In addition, comprehensive evaluation frameworks and tools are developed (Murakonda & Shokri, 2020; Song & Mittal, 2021). Finally, Nasr et al. (2018) proposes a defense mechanism to counter MIAs from an adversarial perspective. Our method shares conceptual similarities with this work, but there are several key differences. Our primary focus is on machine unlearning problems, while their focus is on defending against MIAs. This means that our framework needs to support multiple types of MIAs to provide a comprehensive evaluation of unlearning, including both neural network (NN)-based and non-NN-based attacks. However, their framework only supports NN-based attacks. Furthermore, NN-based attacks are generally not suitable for our runtime requirements; indeed, if unlearning takes longer than retraining, we would opt for retraining instead.

## 3 PRELIMINARIES

**Machine Unlearning.** Let  $D = \{(x_i, y_i) \mid x_i \in \mathcal{X}, y_i \in \mathcal{Y}\}$  be a labeled dataset, where  $\mathcal{X}$  (resp.  $\mathcal{Y}$ ) denote the feature (resp. label) space. The training, validation, and test sets are  $D_{tr}$ ,  $D_{val}$ , and  $D_{te}$ , respectively. A machine learning (ML) algorithm is denoted by  $\mathcal{A}$ , mapping from the joint space of features and labels  $\mathcal{X} \times \mathcal{Y}$  to a hypothesis class. We refer to the model trained on the entire training set as the original model, i.e.,  $\theta_o = \mathcal{A}(D_{tr})$ .

Let  $D_f = \{(x_j^f, y_j^f)\}_{j=1}^q \subseteq D_{tr}$  represent a forget set. The goal of machine unlearning is to remove the influence of  $D_f$  from the original model, resulting in an unlearned model  $\theta_u$  (i.e.,  $\theta_u = \mathcal{U}(\theta_o)$ ) where  $\mathcal{U}$  represents a machine unlearning algorithm. The unlearning algorithm may have access to

other inputs (e.g., the validation set  $D_{val}$ ) depending on the problem settings. Let  $D_r$  be the retain set, the subset of the training data excluding the forget set, i.e.,  $D_r = D_{tr} \setminus D_f$ . The gold standard of machine unlearning is  $\theta_r = \mathcal{A}(D_r)$ , a model trained on the retain set, excluding the influence of  $D_f$ . We denote  $\theta_r$  as the gold standard when comparing machine unlearning algorithms. Retraining is expensive, especially for deep neural networks. This motivates the development of efficient machine unlearning algorithms that satisfy the following conditions: 1) the influence of  $D_f$  is removed from the unlearned model, 2) the performance of the unlearned model is comparable to the performance of the original model, and 3) the computational costs (e.g., running time) are lower compared to those incurred during retraining.

**Membership Inference Attacks.** A membership inference attack (MIA) aims to determine whether a data instance was used to train an ML model (Shokri et al., 2017). An instance that was in the training set is called a member, while one that was not in the training set is called a non-member. Formally, given a target model  $\theta$ , an attacker infers the membership of an instance  $(x, y)$  based on the model’s outputs (i.e.,  $S_\theta(x)$ ) and the label. The attacker does not have access to either the training data or the model parameters of the target model. Instead, he gathers proxy training and test sets and learns a model  $\tilde{\theta}$  to mimic the behavior of the target model. Using the outputs of  $\tilde{\theta}$  on its own training and test data, the attacker acquires a labeled (member v.s. non-member) dataset, and then uses the labeled dataset to train a binary classifier for determining the membership of an instance.

We adapt the idea of MIA to determine whether the influence of the forget set still exists in an unlearned model  $\theta_u$ . Define an auditing set  $\tilde{D}_{\theta_u} = \{(s_j^f, 1), (s_j^{te}, 0)\}_{j=1}^q$ , where  $s_j^f$  (resp.  $s_j^{te}$ ) represents the outputs of the forget (resp. test) instances from the unlearned model, that is,  $s_j^f = S_{\theta_u}(x_j^f)$  (resp.  $s_j^{te} = S_{\theta_u}(x_j^{te})$ ).

Here, the test instances serve as an empirical distribution for the unseen data. The outputs can be scalars, such as the instance-wise cross-entropy losses. The outputs can also be the vectors of probabilities across the classes (Shokri et al., 2017; Carlini et al., 2022). The labels “1” and “0” indicate members and non-members, respectively. The MIA reduces to a binary classification task on  $\tilde{D}_{\theta_u}$ , aiming to differentiate the forget instances from the test ones based on the outputs.

## 4 THE GAME MODEL

We model the machine unlearning problem as a Stackelberg game (SG) between an unlearner who deploys models as services, and an auditor who launches MIAs against the models. The key idea is to assess the effectiveness of an unlearning algorithm by measuring whether the auditor will succeed. In particular, the unlearning is considered effective when the auditor is unable to differentiate between the forget set from the test set based on their outputs from the unlearned model. The SG is played in a sequential manner: the unlearner first deploys an unlearned model, and then the auditor launches an MIA in response. Importantly, the advantage of first-mover endows the unlearner with the power to make a decision knowing that the auditor will play a best response (i.e., launching a strong attack). We now formally define the models for both players.

### 4.1 THE AUDITOR’S MODEL

We begin by defining the auditor’s model. Suppose the unlearner has deployed an unlearned model  $\theta_u$ . Following standard setup (Shokri et al., 2017; Song & Mittal, 2021), we assume that the auditor has black-box access to the model,<sup>1</sup> allowing him to query the model, e.g., submitting data to the model and collecting the outputs. The auditor’s goal is to determine whether the influence of the forget set still exists in the model based on the outputs. To achieve this, the auditor constructs an auditing set  $\tilde{D}_{\theta_u}$ , consisting of the model’s outputs when passing the forget and test instances through the unlearning model  $\theta_u$  (see Section 3 for details about the auditing set). The auditor assesses the distinctiveness of the two sets with a binary classifier trained on the auditing set through cross validation.

Let  $U_a$  be the auditor’s utility function, quantifying the distinctiveness of the forget and test instances. Intuitively, a large  $U_a$  indicates that the outputs of the forget instances are highly differentiable from

<sup>1</sup>

the outputs of the test instances, strong evidence that the influence of the forget set still exists in the unlearned model. We formulate the auditor’s model as the following optimization problem

$$U_a(\theta_a, \theta_u) = M(\tilde{D}_{\theta_u}^{val}; \theta_a) \quad \text{where } \theta_a \in \mathcal{B}_{\theta_u} = \arg \max_{\theta_a \in \mathcal{H}_a} M(\tilde{D}_{\theta_u}^{tr}; \theta_a). \quad (1)$$

The auditing set  $\tilde{D}_{\theta_u}$  is divided into the training  $\tilde{D}_{\theta_u}^{tr}$  and the validation  $\tilde{D}_{\theta_u}^{val}$  sets. The constraint encodes the process of learning a binary classifier. The set  $\mathcal{B}_{\theta_u}$  are the auditor’s best-responses to the unlearner’s decision  $\theta_u$ , that is, a specific MIA that maximally differentiates the forget and test instances.<sup>2</sup> The function  $M$  is an evaluation metric for the binary classifier on a dataset. The definition of  $M$  is flexible. One can use the accuracy to quantify the average performance of the classifier, where true positives are weighted equally with true negatives (Shokri et al., 2017; Song & Mittal, 2021). Alternatively, an average measure may not capture real privacy threats. Instead, ROC curve or true positive rates at specified false positive rates can be used for evaluation Carlini et al. (2022).

The auditor’s model exhibits a high degree of generality, unifying several advanced MIAs in the literature; this includes neural network-based attacks proposed by Nasr et al. (2018), quantile regression-based attacks from Bertram et al. (2019), and prediction confidence-based attacks by Song & Mittal (2021), etc. Under the formulation of equation 1, the mentioned attacks differ in 1) the hypothesis class  $\mathcal{H}_a$  of the binary classifier and 2) the objective function  $M$ . Notice the dependence of the auditor’s best-response on  $\theta_u$  (i.e.,  $\mathcal{B}_{\theta_u}$ ) arising from the unlearner’s first-mover advantage. The unlearner utilizes this dependence to select an unlearned model that limits the auditor’s discriminative power, which we discuss next.

#### 4.2 THE UNLEARNER’S MODEL

Next, we define the unlearner’s model. Let  $C_u$  represent the unlearner’s cost function, which encompasses two main objectives for unlearning. The first objective is to maintain the utility of the model, ensuring that the unlearned model performs comparably (e.g., in terms of predictive power) to the original model. To achieve this objective, we minimize a loss function  $L(D_r; \theta_u)$  computed on the retain set  $D_r$ , following the principles of empirical risk minimization. All regularization terms are included in the loss function to simplify notation. The second objective focuses on eliminating the influence of the forget set from the model being unlearned. We approach this objective adversarially by considering the auditor’s utility  $M$ . In essence, a smaller value of the auditor’s utility indicates that the forget set is harder to be distinguished from the test set, providing strong evidence that the unlearning process is effective.

Formally, the unlearner’s optimization problem is to minimize the cost function below

$$C_u(\theta_u, \theta_a) = L(D_r; \theta_u) + \alpha \cdot M(\tilde{D}_{\theta_u}^{val}; \theta_a). \quad (2)$$

The parameter  $\alpha \in \mathbb{R}^+$  balances the loss  $L$  and the auditor’s utility  $M$ . Depending on the specific setting, the cost function  $C_u$  can be extended to incorporate additional objectives for unlearning. For instance, one can specify that the unlearned model should perform poorly on the forget set (Graves et al., 2021b); this can be achieved by minimizing an evaluation metric (e.g., likelihood) on the forget set. Also, several sparsity-promoting techniques have been shown helpful for unlearning (Jia et al., 2023); one way to achieve this is by adding an  $\ell_1$  regularization to the cost function.

#### 4.3 THE STACKELBERG GAME

Now, with the unlearner and the auditor models in place, we formally define the Stackelberg Game (SG). The SG is to solve the following bi-level optimization problem (Colson et al., 2007)

$$\min_{\theta_u \in \mathcal{H}_u} \underbrace{L(D_r; \theta_u) + \alpha \cdot M(\tilde{D}_{\theta_u}^{val}; \theta_a)}_{\text{Unlearner}} \quad s.t. \quad \underbrace{\theta_a \in \mathcal{B}_{\theta_u}}_{\text{Auditor}}. \quad (3)$$

The objective function has two components: the first ensures generalization by minimizing loss on the retain set, while the second quantifies privacy leakage by assessing the auditor’s ability to differentiate

<sup>2</sup>For example, Bertram et al. (2019) proposed a quantile-regression-based MIA. In this case, the best response is the optimal model parameters for the regression.

between forget and test instances. A lower auditor utility indicates more effective unlearning, as it reduces the distinguishability between the two sets. The hierarchical structure encodes the sequential order of the play, with the upper level corresponding to the unlearner’s optimization problem and the lower level capturing the auditor’s best-responses. During the unlearning, the unlearner needs to proactively consider the auditor’s responses. This requires selecting an unlearning model where the influence of the forget set is erased, or from the auditor’s perspective, the forget instances are indistinguishable from the test ones.

The key assumption of the SG is that if the forget set cannot be distinguished from the test set—in terms of the effectiveness of an MIA—its influence is deemed eliminated from the unlearned model. We justify this assumption from three angles. First, one common way to measure forgetfulness is by assessing the accuracy of the unlearned model on the forget set (Graves et al., 2021b; Chundawat et al., 2023; Baumhauer et al., 2022). This approach is grounded on the observation that machine learning models exhibit distinct performance between training data and unseen data. However, it is important to note that accuracy on the forget set does not necessarily correlate with forgetfulness, as there are inherently difficult (or easy) instances that result in low (or high) accuracy regardless of whether they were part of the training set (Carlini et al., 2022). Secondly, MIAs have been used to study training data forgetting (Jagielski et al., 2023), demonstrating its utility in detecting residual traces of a dataset. Finally, from an adversarial perspective, if a sophisticated attack like an MIA cannot differentiate the forget set from the test set, it is reasonable to expect that the influence of the forget set has been removed.

We solve the SG using gradient-based methods, allowing for easy integration into end-to-end training pipelines. Specifically, we use Implicit Function Theorem to differentiate through the auditor’s optimization problem equation 1, obtaining the gradient of the auditor’s utility with respect to (w.r.t) the unlearning model’s weights, i.e.,  $\partial M / \partial \theta_u$ . As a result, the SG becomes a differentiable layer, compatible with the standard forward-backward computation. The solution methods will be detailed in the next section.

## 5 SG-UNLEARN: STACKELBERG GAME UNLEARN

In this section, we describe the algorithm for solving the SG. In general, it is NP-hard to find an optimal solution for the unlearner (Conitzer & Sandholm, 2006). Instead, we focus on gradient-based algorithms to find an approximate solution, i.e., a model parameter  $\theta_u$  exhibiting good unlearning performance. The main technical challenge is computing the gradient of the auditor’s utility w.r.t. the unlearning model’s weights (i.e.,  $\partial M / \partial \theta_u$ ), which requires differentiation through the auditor’s optimization problem. While the differentiation can be bypassed in some special cases, e.g., when the unlearner’s hypothesis class is of linear regressions (Tong et al., 2018), this is rarely applicable in the current setting given our primary focus on unlearning deep neural networks.

Our solution leverages both the Implicit Function Theorem (IFT) (Dontchev et al., 2009) and tools from Differentiable Optimization (DO) to compute the gradients (Gould et al., 2016; Amos & Kolter, 2017; Agrawal et al., 2019), thereby rendering the SG a differentiable layer seamlessly integrable into existing end-to-end pipelines.

We start by expanding the gradient of  $C_u$  w.r.t.  $\theta_u$  using the chain rule

$$\frac{\partial C_u}{\partial \theta_u} = \frac{\partial L(D_r; \theta_u)}{\partial \theta_u} + \frac{\partial M(\tilde{D}_{\theta_u}^{val}; \theta_a)}{\partial \theta_a} \cdot \frac{\partial \theta_a}{\partial \tilde{D}_{\theta_u}^{tr}} \cdot \frac{\partial \tilde{D}_{\theta_u}^{tr}}{\partial \theta_u}. \quad (4)$$

The first term on the right-hand side can be easily computed using an automatic differentiation tool like PyTorch (Paszke et al., 2017). In essence, the computation involves passing  $D_r$  through the unlearning model (i.e.,  $\theta_u$ ) in the forward pass, computing the loss  $L$ , and getting the gradients in the backward pass. The second term on the right is an expansion of  $\partial M / \partial \theta_u$  using the chain rule; for clarity we omit the arguments of the functions. The gradient  $\partial M / \partial \theta_a$  is obtained by performing a standard forward-backward pass. Some evaluation metrics for binary classification, such as the 0-1 loss, AUC, recall, etc., are non-differentiable. Therefore, we adhere to standard practices by employing a differentiable proxy for  $M$ , such as utilizing the logistic loss as a substitute for the 0-1 loss.

**Leveraging Implicit Function Theorem** Computing the gradient  $\partial\theta_a/\partial\tilde{D}_{\theta_u}^{tr}$  requires differentiation through the attacker’s optimization problem. The main challenge is the absence of an explicit function that maps  $\tilde{D}_{\theta_u}^{tr}$  to  $\theta_a$ . However, under certain regularity assumptions, one can derive an implicit mapping between  $\tilde{D}_{\theta_u}^{tr}$  and  $\theta_a$  based on the optimality conditions of the auditor’s optimization problem (Gould et al., 2016). A concrete example is when the optimization problem is convex, such as learning a Support Vector Machine<sup>3</sup>. In this case, the KKT conditions are necessary and sufficient conditions for the optimality, and it connects  $\theta_a$  with  $\tilde{D}_{\theta_u}^{tr}$  through a system of linear equations, i.e.,

$$f(\tilde{D}_{\theta_u}^{tr}, \theta_a) = \mathbf{0}, \quad (5)$$

where  $f$  encapsulates the stationarity conditions, the primal and dual feasibility conditions, and the complementary slackness conditions (Boyd & Vandenberghe, 2004). For illustration purposes, a concrete example of the KKT conditions  $f$  for linear SVM is provided in Appendix A.9. We apply IFT to the system of linear equations, resulting in

$$\frac{\partial\theta_a}{\partial\tilde{D}_{\theta_u}^{tr}} = - \left( \frac{\partial f(\tilde{D}_{\theta_u}^{tr}, \theta_a)}{\partial\theta_a} \right)^{-1} \frac{\partial f(\tilde{D}_{\theta_u}^{tr}, \theta_a)}{\partial\tilde{D}_{\theta_u}^{tr}}. \quad (6)$$

For further insights into differentiating through an optimization problem using the implicit function theorem, we recommend referring to the lectures by Gould (2023).

**Leveraging Differentiable Optimization** In practice, we capitalize on tools from Differentiable Optimization (DO) to compute the gradients. Intuitively, we can consider DO as software that implements IFT, as shown in equation equation 6, for a given optimization problem. What we need to do is describing the auditor’s optimization problem using a specialized modeling language, e.g., `cvxpy` (Diamond & Boyd, 2016). We then use DO to transform this description into a differentiable layer. Subsequently, this differentiable layer is positioned atop the model undergoing unlearning, thereby establishing a computational pathway from  $\theta_u$  to  $\theta_a$ . The pseudocode for this process is provided in Algorithm 1. This algorithm has a time complexity of  $O(n^3)$ , where  $n$  denotes the size of the attacker’s optimization problem (i.e., the number of variables and/or constraints). This cubic dependence stems from the matrix inversion in equation 6.

---

**Algorithm 1** SG-Unlearn

---

- 1: Input:  $D_r, D_f, D_{te}$  and the original model  $\theta_o$
  - 2: Initialize:  $i = 0, \theta_u^0 = \theta_o$ , a scheduler  $\eta^i$
  - 3: **while**  $i < \text{epoch}$  **do**
  - 4:   Compute  $L(D_r; \theta_u^i)$  on the retain set in a forward pass
  - 5:   Update  $\theta_u^{i'} \leftarrow \theta_u^i - \eta^i \cdot \frac{\partial L(D_r; \theta_u^i)}{\partial\theta_u^i}$
  - 6:   Construct the auditing set  $\tilde{D}_{\theta_u^{i'}}$  from  $D_f$  and  $D_{te}$
  - 7:   Describe the auditor’s optimization problem equation 1 with `cvxpy`
  - 8:   Convert the description to a differentiable layer `AuditorLayer`
  - 9:   Plug `AuditorLayer` into the computational graph
  - 10:   Get the auditor’s best response  $\theta_a^i \leftarrow \text{AuditorLayer}(\tilde{D}_{\theta_u^{i'}})$
  - 11:   Compute  $M(\tilde{D}_{\theta_u^{i'}}^{val}; \theta_a^i)$
  - 12:   Update  $\theta_u^{i+1} \leftarrow \theta_u^{i'} - \eta^i \cdot \frac{\partial M(\tilde{D}_{\theta_u^{i'}}^{val}; \theta_a^i)}{\partial\theta_u^{i'}}$
  - 13:    $i \leftarrow i + 1$
  - 14: **end while**
  - 15: Return:  $\theta_u^i$
- 

## 6 EXPERIMENTS

### 6.1 EXPERIMENT SETUP

We conduct experiments on both computer vision (CV) and natural language processing (NLP) datasets. For the CV tasks, we use the widely recognized image classification datasets CIFAR-10,

<sup>3</sup>This includes several state-of-the-art MIAs (Bertran et al., 2023; Song & Mittal, 2021).

CIFAR-100, and SVHN (Krizhevsky et al., 2009; Netzer et al., 2011). The backbone model we use is ResNet-18 (He et al., 2016). For NLP tasks, we assess performance on the 20 Newsgroups dataset, leveraging the BERT model. We explore two learning scenarios: *random forgetting* and *class-wise forgetting*. In random forgetting, instances are sampled uniformly at random from all classes. In contrast, class-wise forgetting involves selecting all instances from a specific class. For CIFAR-10 and CIFAR-100, the forget set consists of 10% of the entire training set, while the ratio is reduced to 5% for SVHN. In all experiments, the attacker’s optimization problem is formulated as a binary classification task, where a linear support vector machine (SVM) is used to distinguish between forget and test instances.

The baseline methods we use for comparison with SG include Retrain<sup>4</sup>, Fine-Tune (FT) (Warnecke et al., 2021; Golatkar et al., 2020a), Gradient Ascent (GA) (Graves et al., 2021a; Thudi et al., 2022a), Influence Unlearning (IU) (Izzo et al., 2021a; Koh & Liang, 2017b),  $\ell_1$ -sparse (Jia et al., 2023), Random Label (RL) (Hayase et al., 2020), Boundary Expansion (BE), Boundary Shrink (BS) (Chen et al., 2023) and SCRUB (Kurmanji et al., 2024). Further details on the baseline methods are provided in Section A.6 of the Appendix. For all methods, we use the SGD optimizer with a weight decay of  $5e-4$  and a momentum of 0.9. Other hyper-parameters are selected through the validation set. Specifically, we create a new auditing set. For each unlearning method, we select the hyper-parameters that maximize the difference between the validation accuracy and the MIA accuracy on this new auditing set. This approach ensures that the model both generalizes well to unseen data (high validation accuracy) and is less vulnerable to the attacks (low MIA accuracy). The hyperparameters are listed in Table 9 in the Appendix.

## 6.2 EVALUATION METRICS

We evaluate SG and the baseline methods using metrics commonly adopted in prior studies (Bourtoule et al., 2021; Jagielski et al., 2023; Jia et al., 2023; Chundawat et al., 2023). *It is important to note that the test accuracy is evaluated on a subset of the test data that is separate from the one used for solving SG.* **Retain accuracy ( $Acc_r$ ) and test accuracy ( $Acc_{te}$ )** are used to quantify model utility (Jia et al., 2023). **MIA accuracy, AUC and F1 score** are the metrics to quantify the effectiveness of unlearning, all of which are estimated on the auditing set with 10-fold cross Carlini et al. (2022). An effective unlearning algorithm should result in MIA metrics that approach random guessing (0.5). **Forget accuracy ( $Acc_f$ ) and the absolute difference between the forget and test accuracy ( $|Acc_f - Acc_{te}|$ )**: An effective unlearning algorithm should result in  $Acc_f$  being close to  $Acc_{te}$ . This indicates that the unlearned model no longer retains specific information about the forget data, as its performance on the forget set should be similar to its performance on unseen test data ( $D_{te}$ ), reflecting the removal of the influence of  $D_f$ . To gather additional statistical evidence regarding the effectiveness of unlearning, we collect the cross-entropy losses of the forget and test instances from the unlearned model into the empirical distributions  $\mathcal{L}_f$  and  $\mathcal{L}_{te}$ , respectively. Next, we run a **Kolmogorov-Smirnov statistics (KS Stat.)** test to determine if the distributions can be differentiated from each other. The KS statistic quantifies the differences between  $\mathcal{L}_f$  and  $\mathcal{L}_{te}$ , where the p-value indicates whether the difference is significant (Massey Jr, 1951). In addition to the KS statistics, we provide the **Wasserstein distance (W. Dist.)** between the empirical distributions of  $\mathcal{L}_f$  and  $\mathcal{L}_{te}$ . This complements the KS statistics and evaluates the unlearning performance in terms of the similarity between the losses.

## 6.3 RESULTS

The experimental results for random forgetting and class-wise forgetting are presented in Section 6.3.1 and 6.3.2, respectively. We consider retrain as the gold standard for evaluating unlearning algorithms: the closer to the metrics of retrain the more effective the algorithm. We highlight the closest metrics to retrain in bold.

### 6.3.1 RANDOM FORGETTING

We present the results for CIFAR-10, CIFAR-100 and 20 NewsGroup in Table 1. The results for SVHN datasets are provided in Appendix A.2.1. SG achieves the best performance for most of

<sup>4</sup>Retraining the unlearning model on  $D_{tr} \setminus D_f$  from scratch.

the metrics, demonstrating its effectiveness in unlearning. Specifically, the KS statistic of SG is consistently lower than those of the other baselines, exhibiting an order of magnitude difference in the statistics for CIFAR-10 compared to most baselines. Intuitively, ML models behave differently on training data compared to unseen data, and this difference is usually reflected in the corresponding losses (Carlini et al., 2022). The small KS statistic of SG implies that the forget and test instances exhibit greater similarity in terms of the model’s behavior, although there is still a discernible difference between the losses. Another metric for measuring the similarity is the Wasserstein distance (W. Dist.). The baseline RL achieves the lowest distance, although the difference with SG is not statistically significant. A visualization of the cross-entropy losses for the forget and test instances from one of the experiments is provided in Figure 4 in the Appendix. The experiment results of CelebA and TinyImageNet are given in Table 7 and 6 in Appendix.

Table 1: Experimental results (Mean<sub>std</sub>) on CIFAR-10, CIFAR-100 and 20 NewsGroup for random forgetting. The highlighted metrics are the closest to those of retraining, which is considered as the best performance compared with the other baselines.

CIFAR-10	$Acc_r$	$Acc_{te}$	$Acc_f$	$ Acc_f - Acc_{te} $	MIA acc.	MIA AUC	MIA F1	KS Stat.	W. Dist.	RTE (min., ↓)
Retrain	0.9996 <sub>0.0001</sub>	0.9291 <sub>0.0022</sub>	0.9230 <sub>0.0043</sub>	0.0061	0.5069 <sub>0.0073</sub>	0.5083 <sub>0.0099</sub>	0.5094 <sub>0.0187</sub>	0.0255 <sub>0.0080</sub>	0.0307 <sub>0.0115</sub>	14.92
FT	0.9886 <sub>0.0055</sub>	0.9114 <sub>0.0050</sub>	0.9851 <sub>0.0056</sub>	0.0737	0.5405 <sub>0.0031</sub>	0.5457 <sub>0.0070</sub>	0.6293 <sub>0.0127</sub>	0.0933 <sub>0.0050</sub>	0.3158 <sub>0.0142</sub>	0.45
GA	<b>0.9996</b> <sub>0.0001</sub>	0.9304 <sub>0.0007</sub>	0.9995 <sub>0.0003</sub>	0.0691	0.5504 <sub>0.0066</sub>	0.5611 <sub>0.0071</sub>	0.6625 <sub>0.0106</sub>	0.1403 <sub>0.0037</sub>	0.2782 <sub>0.0026</sub>	0.18
IU	0.9723 <sub>0.0255</sub>	0.8966 <sub>0.0242</sub>	0.9722 <sub>0.0243</sub>	0.0756	0.5398 <sub>0.0055</sub>	0.5548 <sub>0.0076</sub>	0.6193 <sub>0.0204</sub>	0.1050 <sub>0.0192</sub>	0.4083 <sub>0.0554</sub>	<b>0.02</b>
$\ell_1$ -sparse	0.9970 <sub>0.0007</sub>	0.9234 <sub>0.0014</sub>	0.9938 <sub>0.0016</sub>	0.0704	0.5501 <sub>0.0049</sub>	0.5694 <sub>0.0087</sub>	0.6405 <sub>0.0259</sub>	0.1018 <sub>0.0034</sub>	0.2704 <sub>0.0074</sub>	0.96
RL	0.9989 <sub>0.0001</sub>	0.9217 <sub>0.0008</sub>	0.9810 <sub>0.0025</sub>	0.0593	0.5217 <sub>0.0087</sub>	0.5297 <sub>0.0133</sub>	<b>0.5935</b> <sub>0.0140</sub>	0.0986 <sub>0.0136</sub>	<b>0.1520</b> <sub>0.0058</sub>	0.84
BE	<b>0.9996</b> <sub>0.0001</sub>	0.9304 <sub>0.0007</sub>	0.9996 <sub>0.0003</sub>	0.0692	0.5541 <sub>0.0049</sub>	0.5639 <sub>0.0058</sub>	0.6629 <sub>0.0082</sub>	0.1412 <sub>0.0030</sub>	0.2783 <sub>0.0020</sub>	0.27
BS	0.9995 <sub>0.0001</sub>	0.9307 <sub>0.0008</sub>	0.9995 <sub>0.0004</sub>	0.0688	0.5588 <sub>0.0072</sub>	0.5779 <sub>0.0097</sub>	0.6590 <sub>0.0156</sub>	0.1466 <sub>0.0032</sub>	0.3072 <sub>0.0026</sub>	0.46
SCRUB	0.9971 <sub>0.0018</sub>	<b>0.9251</b> <sub>0.0018</sub>	0.9959 <sub>0.0022</sub>	0.0708	0.5533 <sub>0.0059</sub>	0.5679 <sub>0.0073</sub>	0.6337 <sub>0.0149</sub>	0.1038 <sub>0.0071</sub>	0.2485 <sub>0.0154</sub>	1.30
GAU	0.9583 <sub>0.0033</sub>	<b>0.9009</b> <sub>0.0026</sub>	0.9311 <sub>0.0042</sub>	0.0302	0.5184 <sub>0.0018</sub>	0.5146 <sub>0.0090</sub>	0.6439 <sub>0.0008</sub>	0.0393 <sub>0.0057</sub>	0.1216 <sub>0.0091</sub>	8.34
SG	0.9948 <sub>0.0029</sub>	0.8940 <sub>0.0048</sub>	0.9351 <sub>0.0070</sub>	0.0411	0.5202 <sub>0.0054</sub>	0.5134 <sub>0.0084</sub>	0.6480 <sub>0.0043</sub>	0.0482 <sub>0.0082</sub>	0.1555 <sub>0.0194</sub>	1.47
SG (Acc.)	0.9962 <sub>0.0003</sub>	0.8870 <sub>0.0011</sub>	0.9090 <sub>0.0001</sub>	0.0293	<b>0.5110</b> <sub>0.0003</sub>	0.5200 <sub>0.0001</sub>	0.6358 <sub>0.0002</sub>	0.0274 <sub>0.0005</sub>	0.1087 <sub>0.0016</sub>	0.88
SG + LiRA	0.9948 <sub>0.0038</sub>	0.8865 <sub>0.0059</sub>	<b>0.9158</b> <sub>0.0093</sub>	0.0293	0.5151 <sub>0.0039</sub>	<b>0.5100</b> <sub>0.0070</sub>	0.6390 <sub>0.0026</sub>	<b>0.0363</b> <sub>0.0059</sub>	0.1126 <sub>0.0014</sub>	8.33
SG (Acc.) + RL	0.9968 <sub>0.0048</sub>	0.9237 <sub>0.0051</sub>	0.9468 <sub>0.0108</sub>	<b>0.0231</b>	0.5208 <sub>0.0145</sub>	0.5230 <sub>0.0206</sub>	<b>0.5236</b> <sub>0.0198</sub>	0.0958 <sub>0.0224</sub>	<b>0.1082</b> <sub>0.0191</sub>	1.79
CIFAR-100	$Acc_r$	$Acc_{te}$	$Acc_f$	$ Acc_f - Acc_{te} $	MIA acc.	MIA AUC	MIA F1	KS Stat.	W. Dist.	RTE (min., ↓)
Retrain	0.9996 <sub>0.0001</sub>	0.7035 <sub>0.0025</sub>	0.6925 <sub>0.0039</sub>	0.0110	0.5184 <sub>0.0057</sub>	0.5281 <sub>0.0053</sub>	0.5104 <sub>0.0081</sub>	0.0203 <sub>0.0045</sub>	0.0567 <sub>0.0200</sub>	13.08
FT	0.9991 <sub>0.0001</sub>	0.7117 <sub>0.0021</sub>	0.9984 <sub>0.0006</sub>	0.2867	0.6630 <sub>0.0075</sub>	0.7300 <sub>0.0102</sub>	0.6878 <sub>0.0107</sub>	0.4566 <sub>0.0083</sub>	1.2583 <sub>0.0166</sub>	0.39
GA	<b>0.9996</b> <sub>0.0001</sub>	0.7158 <sub>0.0008</sub>	0.9996 <sub>0.0002</sub>	0.2838	0.6977 <sub>0.0060</sub>	0.7601 <sub>0.0065</sub>	0.7207 <sub>0.0088</sub>	0.4915 <sub>0.0030</sub>	1.2219 <sub>0.0038</sub>	<b>0.20</b>
IU	0.9971 <sub>0.0029</sub>	<b>0.7026</b> <sub>0.0080</sub>	0.9959 <sub>0.0034</sub>	0.2933	0.6660 <sub>0.0089</sub>	0.7305 <sub>0.0134</sub>	0.6950 <sub>0.0124</sub>	0.4583 <sub>0.0168</sub>	1.2612 <sub>0.0366</sub>	0.21
$\ell_1$ -sparse	0.9958 <sub>0.0013</sub>	0.7095 <sub>0.0025</sub>	0.9890 <sub>0.0028</sub>	0.2785	0.6738 <sub>0.0081</sub>	0.7392 <sub>0.0088</sub>	0.6952 <sub>0.0073</sub>	0.3717 <sub>0.0113</sub>	1.1157 <sub>0.0079</sub>	0.84
RL	0.9965 <sub>0.0054</sub>	0.6665 <sub>0.0031</sub>	0.8483 <sub>0.0447</sub>	0.1818	0.5808 <sub>0.0426</sub>	0.6152 <sub>0.0527</sub>	0.6080 <sub>0.0466</sub>	0.2323 <sub>0.0731</sub>	0.8067 <sub>0.1705</sub>	0.73
BE	0.9995 <sub>0.0001</sub>	0.7173 <sub>0.0014</sub>	0.9996 <sub>0.0002</sub>	0.2823	0.6977 <sub>0.0037</sub>	0.7661 <sub>0.0066</sub>	0.7248 <sub>0.0065</sub>	0.4940 <sub>0.0031</sub>	1.2175 <sub>0.0119</sub>	0.23
BS	0.9995 <sub>0.0001</sub>	0.7160 <sub>0.0013</sub>	0.9996 <sub>0.0002</sub>	0.2836	0.6987 <sub>0.0052</sub>	0.7651 <sub>0.0072</sub>	0.7239 <sub>0.0054</sub>	0.4963 <sub>0.0033</sub>	1.2382 <sub>0.0206</sub>	0.39
SCRUB	0.9993 <sub>0.0001</sub>	0.7097 <sub>0.0019</sub>	0.9991 <sub>0.0004</sub>	0.2894	0.7015 <sub>0.0057</sub>	0.7747 <sub>0.0060</sub>	0.7299 <sub>0.0056</sub>	0.4717 <sub>0.0052</sub>	1.2280 <sub>0.0128</sub>	1.14
GAU	0.9346 <sub>0.0006</sub>	0.6687 <sub>0.0045</sub>	0.8021 <sub>0.0073</sub>	0.1334	0.5760 <sub>0.0044</sub>	0.5816 <sub>0.0045</sub>	0.6526 <sub>0.0042</sub>	0.1380 <sub>0.0134</sub>	0.7268 <sub>0.0858</sub>	13.35
SG	0.8993 <sub>0.0105</sub>	0.6378 <sub>0.0066</sub>	0.7239 <sub>0.0093</sub>	0.0861	0.5412 <sub>0.0070</sub>	0.5320 <sub>0.0076</sub>	0.6061 <sub>0.0056</sub>	0.0988 <sub>0.0061</sub>	0.5316 <sub>0.0295</sub>	3.07
SG (Acc.)	0.9646 <sub>0.0019</sub>	0.6028 <sub>0.0003</sub>	0.7032 <sub>0.0027</sub>	0.1004	0.5519 <sub>0.0008</sub>	0.5571 <sub>0.0010</sub>	0.6192 <sub>0.0004</sub>	0.1120 <sub>0.0029</sub>	0.6176 <sub>0.0007</sub>	1.55
SG + LiRA	0.9574 <sub>0.0038</sub>	0.6008 <sub>0.0059</sub>	<b>0.6942</b> <sub>0.0093</sub>	<b>0.0293</b>	0.5411 <sub>0.0039</sub>	<b>0.5303</b> <sub>0.0070</sub>	0.6057 <sub>0.0026</sub>	<b>0.0274</b> <sub>0.0059</sub>	<b>0.1087</b> <sub>0.0014</sub>	10.68
SG (Acc.) + RL	0.9966 <sub>0.0014</sub>	0.6679 <sub>0.0033</sub>	0.8095 <sub>0.0113</sub>	0.1416	0.5609 <sub>0.0397</sub>	0.6032 <sub>0.0412</sub>	<b>0.6013</b> <sub>0.0023</sub>	0.2041 <sub>0.0935</sub>	0.6391 <sub>0.1845</sub>	2.39
20 NewsGroup	$Acc_r$	$Acc_{te}$	$Acc_f$	$ Acc_f - Acc_{te} $	MIA acc.	MIA AUC	MIA F1	KS Stat.	W. Dist.	RTE (min., ↓)
Retrain	1.0000	0.8528	0.9224	0.0696	0.5285	0.5512	0.5501	0.1405	0.5925	40.8
FT	0.9999	0.8518	0.8035	<b>0.0482</b>	0.5672	0.6059	0.6220	0.2495	1.1894	20.2
GA	0.0483	0.0483	0.0500	0.0017	0.4995	<b>0.4973</b>	0.2704	0.0334	0.0990	26.1
IU	<b>1.0000</b>	<b>0.8575</b>	<b>0.9990</b>	0.1415	0.5676	0.6054	0.6348	0.2986	<b>0.9614</b>	27.9
RL	0.9985	0.8298	0.6709	0.1589	0.7123	0.7651	0.7148	0.5334	1.1402	21.2
SG	<b>1.0000</b>	1.0000	1.0000	0.0000	<b>0.5065</b>	0.4922	<b>0.5627</b>	<b>0.0791</b>	0.0007	<b>15.6</b>

Another observation from the table is the inherent trade-off between model performance, measured by test accuracy, and the effectiveness of unlearning, measured by MIA accuracy. This trade-off has been documented in prior studies as a common challenge in unlearning tasks (Golatkar et al., 2020a; Bourtole et al., 2021). Specifically, SG is more effective at unlearning the forget instances, as indicated by the highlighted MIA metrics. However, this effectiveness comes at a cost to the test accuracy on CIFAR-10 and CIFAR-100, a phenomenon observed in other unlearning techniques as well (Jia et al., 2023; Graves et al., 2021a). Despite this trade-off, the degradation in test accuracy remains minimal. We run a large array of experiments with varying  $\alpha$  from  $\{0.05, 0.1, 0.25, 0.5, 1, 2, 5\}$  (see equation 2) to explore the extent to which the trade-off can be reduced. The results are presented in Figure 2. Unfortunately, we do not see a consistent trend that makes SG closer to retrain across all the metrics.

### 6.3.2 CLASS-WISE FORGETTING

We use CIFAR-10 as the benchmark for class-wise forgetting. For results on other datasets, please refer to Appendix A.8. In random forgetting, instances are uniformly sampled across all classes, preserving the overall dataset distribution. In contrast, class-wise forgetting removes all instances from a specific class, resulting in a significant distribution shift that makes forget data more detectable. The experimental results, presented in Table 2, highlight the metrics closest to retraining. However, all



methods perform poorly on MIA-related metrics, as the auditor can easily distinguish between forget and test instances due to the distinct distributional shifts caused by class-wise forgetting. Additionally, no single method consistently outperforms the others.

Table 2: Experimental results (Mean<sub>std</sub>) on CIFAR-10 for class-wise forgetting. The highlighted metrics are the closest to those of retrain, which is considered as the best performance compared with the other baselines.

CIFAR-10	$Acc_r$	$Acc_{te}$	$Acc_f$	$ Acc_f - Acc_{te} $	MIA acc.	MIA AUC	MIA F1	KS Stat.	W. Dist.	RTE (min., ↓)
Retrain	0.9996 <sub>0.0001</sub>	0.9333 <sub>0.0009</sub>	0.0000 <sub>0.0000</sub>	0.9333	0.9935 <sub>0.0006</sub>	0.9983 <sub>0.0004</sub>	0.9936 <sub>0.0007</sub>	0.9803 <sub>0.0002</sub>	9.5601 <sub>0.0911</sub>	13.96
FT	0.9958 <sub>0.0022</sub>	0.9226 <sub>0.0030</sub>	0.6043 <sub>0.0450</sub>	0.3183	0.9915 <sub>0.0011</sub>	<b>0.9985</b> <sub>0.0002</sub>	0.9915 <sub>0.0011</sub>	0.7975 <sub>0.0133</sub>	0.9324 <sub>0.1847</sub>	1.16
GA	0.8478 <sub>0.0046</sub>	0.7942 <sub>0.0055</sub>	0.0007 <sub>0.0002</sub>	0.7935	<b>0.9944</b> <sub>0.0011</sub>	0.9996 <sub>0.0002</sub>	<b>0.9938</b> <sub>0.0015</sub>	0.9269 <sub>0.0087</sub>	15.2941 <sub>0.1656</sub>	0.84
IU	0.9339 <sub>0.0161</sub>	0.8644 <sub>0.0141</sub>	0.0619 <sub>0.0149</sub>	0.8025	0.9972 <sub>0.0009</sub>	0.9996 <sub>0.0001</sub>	0.9972 <sub>0.0007</sub>	0.8151 <sub>0.0195</sub>	<b>8.2574</b> <sub>0.6484</sub>	<b>0.31</b>
$\ell_1$ -sparse	0.9972 <sub>0.0005</sub>	0.9285 <sub>0.0014</sub>	0.0914 <sub>0.0310</sub>	0.8371	0.9910 <sub>0.0014</sub>	0.9989 <sub>0.0001</sub>	0.9910 <sub>0.0014</sub>	0.9208 <sub>0.0078</sub>	2.5552 <sub>0.1738</sub>	1.84
RL	<b>0.9996</b> <sub>0.0000</sub>	<b>0.9330</b> <sub>0.0008</sub>	0.0001 <sub>0.0001</sub>	<b>0.9329</b>	0.9916 <sub>0.0013</sub>	0.9990 <sub>0.0005</sub>	0.9916 <sub>0.0013</sub>	0.9695 <sub>0.0025</sub>	6.3989 <sub>0.0789</sub>	1.97
BE	0.9710 <sub>0.0012</sub>	0.8984 <sub>0.0023</sub>	0.2477 <sub>0.0022</sub>	0.6507	0.9964 <sub>0.0005</sub>	0.9990 <sub>0.0004</sub>	0.9964 <sub>0.0005</sub>	0.7306 <sub>0.0047</sub>	4.8984 <sub>0.0432</sub>	0.32
BS	0.9691 <sub>0.0031</sub>	0.8969 <sub>0.0031</sub>	0.2504 <sub>0.0105</sub>	0.6465	0.9965 <sub>0.0001</sub>	0.9988 <sub>0.0005</sub>	0.9965 <sub>0.0002</sub>	0.7196 <sub>0.0072</sub>	5.0155 <sub>0.0922</sub>	0.66
SCRUB	1.0000 <sub>0.0000</sub>	0.9269 <sub>0.0021</sub>	<b>0.0000</b> <sub>0.0000</sub>	0.9269	1.0000 <sub>0.0000</sub>	1.0000 <sub>0.0000</sub>	1.0000 <sub>0.0000</sub>	0.9999 <sub>0.0001</sub>	70.9934 <sub>2.9441</sub>	3.47
SG	0.9667 <sub>0.0054</sub>	0.9056 <sub>0.0055</sub>	<b>0.0000</b> <sub>0.0000</sub>	0.9056	0.9814 <sub>0.0026</sub>	0.9902 <sub>0.0025</sub>	0.9818 <sub>0.0025</sub>	<b>0.9696</b> <sub>0.0032</sub>	5.2754 <sub>0.1882</sub>	0.84

### 6.3.3 THE EFFECT OF THE ATTACKER MODEL

Finally, we conduct a comparative study to understand the impact of adversarial modeling on the unlearning process, controlled by the parameter  $\alpha$  as defined in equation 2. We show the results for random forgetting and defer the results for class-wise forgetting to the appendix. In Figure 1, we compare two cases where  $\alpha$  is set to either 1 or 0, denoted by SG-1 and SG-0 respectively. The comparison is done across four metrics: 1) the test accuracy; 2) the MIA accuracy; 3) the defender’s utility, evaluated as the test accuracy minus the MIA accuracy, which provides a combined scalar value that measures both the performance of the unlearned model and the effectiveness of unlearning; 4) the Wasserstein distance between the empirical distributions of  $\mathcal{L}_f$  and  $\mathcal{L}_{te}$ . We show the averages over 10 experiments with different seeds, and 95% confidence intervals are displayed. The first observation is that the adversarial term (i.e.,  $\alpha \cdot M(\tilde{D}_{\theta_u}^{val}; \theta_a)$ ) acts as a regularizer, improving the generalizability of the unlearned model. This observation is supported by comparing the test accuracy of SG-1 and SG-0 on CIFAR-10 (top middle). Similar findings have been reported in Nasr et al. (2018). Another observation is that adversarial modeling limits the attacker’s ability to differentiate between forget instances and test instances; this is demonstrated by the MIA accuracy on CIFAR-100. The right-most column displays the Wasserstein distances between  $\mathcal{L}_f$  and  $\mathcal{L}_{te}$ . It is evident that the two losses are closer as a result of adversarial modeling, especially for CIFAR-100 dataset. Additionally, the distances progressively decrease throughout the epochs, confirming the effectiveness of the gradient-based method.

In addition to the existence of attacker, we also investigate the strength of the attacker by changing  $\alpha$ . We select the  $\alpha$  in large range of  $\{0.05, 0.1, 0.25, 0.5, 1, 2, 5\}$ . In Figure 2, we compare the performance regarding the test accuracy  $Acc_{te}$ . The cross of the red dash line is the performance of the retrain model. We can find that SG is robust to the attacker strength.

## 6.4 MIA SELECTION

To validate the generalization ability of SG, we select the MIA used in (Jia et al., 2023) as the attacker and we evaluate SG according to (Jia et al., 2023). The results given in Table 3 show that SG is not sensitive to the MIA.

## 7 DISCUSSION

In this paper, we design an adversarial framework for addressing the problem of unlearning a set data from a machine learning model. Our approach focuses on evaluating the effectiveness of unlearning from an adversarial perspective, leveraging membership inference attacks (MIAs) to detect any residual traces of the data within the model. The framework allows for a proactive design of the unlearning algorithm, synthesizing two lines of research—machine unlearning and MIAs—that have heretofore progressed in parallel. By using implicit differentiation techniques, we develop a gradient-based algorithm for solving the game, making the framework easily integrable into existing end-to-end learning pipelines. We present empirical results to support the efficacy of the framework and the algorithm. We believe our work can make a progress in trustworthy ML.

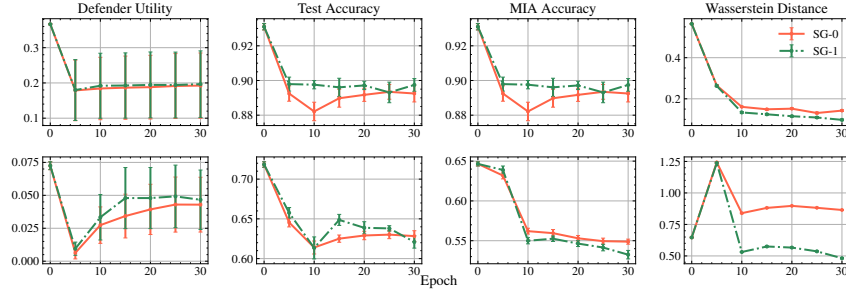


Figure 1: An ablation study to understand the impact of adversarial modeling on the process of unlearning;  $\alpha = 1$  and  $\alpha = 0$  corresponds to the cases with and without adversarial modeling, respectively. The results are the averages over 10 experiments with different seeds, and 95% confidence intervals are displayed. **From the left to the right:** 1) the defender’s utility, evaluated as the test accuracy  $Acc_{te}$  minus the MIA accuracy; 2) test accuracy; 3) MIA accuracy; 4) Wasserstein distance between the cross-entropy losses of the forget and test instances. **Top row:** CIFAR-10; **Bottom row:** CIFAR-100. **Epoch 0: Original model.**

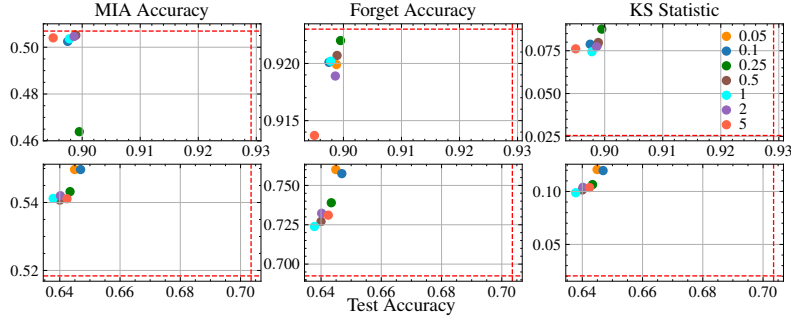


Figure 2: Experiments with different values of the trade-off parameter  $\alpha$ . We consider 7 values  $\{0.05, 0.1, 0.25, 0.5, 1, 2, 5\}$ . Each dot represents a batch of 5 random experiments with the same  $\alpha$ . The coordinates of a dot are the corresponding metrics averaged over the 5 runs. **Top row:** CIFAR-10; **Bottom row:** CIFAR-100.

Table 3: SG is evaluated using the MIA attacker described in (Jia et al., 2023), comparing its performance against baseline models on the CIFAR-10 dataset under the random forgetting paradigm. For the metrics UA, MIA, RA, and TA, the value closest to the Retrain baseline is highlighted in bold.

CIFAR-10	UA ( $1 - Acc_f$ )	MIA	RA ( $Acc_r$ )	TA ( $Acc_r$ )	Avg. Gap ( $\downarrow$ )	RTE (min, $\downarrow$ )
Retrain	0.0807 <sub>0.0047</sub>	0.1741 <sub>0.0069</sub>	1.0000 <sub>0.0001</sub>	0.9161 <sub>0.0024</sub>	-	24.66
FT	0.0110 <sub>0.0019</sub>	0.0406 <sub>0.0041</sub>	0.9983 <sub>0.0003</sub>	0.9370 <sub>0.0010</sub>	0.0555	1.58
GA	0.0056 <sub>0.0001</sub>	0.0119 <sub>0.0005</sub>	0.9948 <sub>0.0002</sub>	0.9455 <sub>0.0005</sub>	0.0680	<b>0.31</b>
IU	0.1751 <sub>0.0219</sub>	0.2139 <sub>0.0170</sub>	0.8328 <sub>0.0244</sub>	0.7813 <sub>0.0285</sub>	0.1091	1.18
$\ell_1$ -sparse	0.0121 <sub>0.0038</sub>	0.0433 <sub>0.0052</sub>	0.9739 <sub>0.0031</sub>	0.9549 <sub>0.0018</sub>	0.0661	1.82
RL	0.0280 <sub>0.0037</sub>	0.1859 <sub>0.0348</sub>	<b>0.9997</b> <sub>0.0001</sub>	0.9408 <sub>0.0012</sub>	0.0224	1.98
BE	0.0000 <sub>0.0000</sub>	0.0026 <sub>0.0002</sub>	1.0000 <sub>0.0000</sub>	0.9535 <sub>0.0018</sub>	0.0724	3.17
BS	0.0048 <sub>0.0007</sub>	0.0116 <sub>0.0004</sub>	0.9947 <sub>0.0001</sub>	0.9458 <sub>0.0003</sub>	0.0684	1.41
SCRUB	0.0070 <sub>0.0059</sub>	0.0388 <sub>0.0125</sub>	0.9959 <sub>0.0034</sub>	0.9422 <sub>0.0026</sub>	0.0598	4.05
SG	<b>0.0748</b> <sub>0.0041</sub>	<b>0.1835</b> <sub>0.0117</sub>	0.9990 <sub>0.0131</sub>	<b>0.9072</b> <sub>0.0189</sub>	<b>0.0063</b>	3.48

**Limitation** One future direction is to enhance the algorithm’s efficiency. As shown in equation 6, the gradient-based algorithm requires a matrix inversion, which exhibits an  $O(n^3)$  dependence on the size of the auditor’s optimization problem. Therefore, developing a more efficient method to differentiate through the auditor’s optimization could significantly accelerate the algorithm. Another direction is to experiment with different combinations of the unlearning algorithm and the MIA within the SG framework. Currently, the unlearner employs Fine-Tune as the unlearning algorithm, while the auditor uses an SVM-based MIA. Exploring the performance of other combinations, such as Random Label with a neural network-based MIA, would be worthwhile.

## ETHIC STATEMENT

This work does not involve potential malicious or unintended uses, fairness considerations, privacy considerations, security considerations, crowd sourcing, or research with human subjects.

## REPRODUCIBILITY STATEMENT

We provide details to reproduce our results in Appendix A.6 and A.7. We also provide pseudo-code in Algorithm 1 and will release the code upon acceptance.

## REFERENCES

- Akshay Agrawal, Brandon Amos, Shane T. Barratt, Stephen P. Boyd, Steven Diamond, and J. Zico Kolter. Differentiable convex optimization layers. In *Proceedings of the 32nd Conference on Neural Information Processing Systems (NeurIPS)*, pp. 9558–9570, 2019.
- Brandon Amos and J Zico Kolter. Optnet: Differentiable optimization as a layer in neural networks. In *Proceedings of the 34th International Conference on Machine Learning (ICML)*, pp. 136–145. PMLR, 2017.
- Thomas Baumhauer, Pascal Schöttle, and Matthias Zeppelzauer. Machine unlearning: Linear filtration for logit-based classifiers. *Machine Learning*, 111(9):3203–3226, 2022.
- Theo Bertram, Elie Bursztein, Stephanie Caro, Hubert Chao, Rutledge Chin Feman, Peter Fleischer, Albin Gustafsson, Jess Hemerly, Chris Hibbert, Luca Invernizzi, et al. Five years of the right to be forgotten. In *Proceedings of the ACM Conference on Computer and Communications Security (CCS)*, pp. 959–972, 2019.
- Martin Bertran, Shuai Tang, Michael Kearns, Jamie Morgenstern, Aaron Roth, and Zhiwei Steven Wu. Scalable membership inference attacks via quantile regression. *CoRR*, abs/2307.03694, 2023. URL <https://doi.org/10.48550/arXiv.2307.03694>.
- Lucas Bourtole, Varun Chandrasekaran, Christopher A Choquette-Choo, Hengrui Jia, Adelin Travers, Baiwu Zhang, David Lie, and Nicolas Papernot. Machine unlearning. In *Proceedings of IEEE Symposium on Security and Privacy (S&P)*, pp. 141–159. IEEE, 2021.
- Stephen P Boyd and Lieven Vandenbergh. *Convex optimization*. Cambridge university press, 2004.
- Nicholas Carlini, Steve Chien, Milad Nasr, Shuang Song, Andreas Terzis, and Florian Tramer. Membership inference attacks from first principles. In *Proceedings of the 39th IEEE Symposium on Security and Privacy (S&P)*, pp. 1897–1914. IEEE, 2022.
- Min Chen, Zhikun Zhang, Tianhao Wang, Michael Backes, Mathias Humbert, and Yang Zhang. When machine unlearning jeopardizes privacy. In Yongdae Kim, Jong Kim, Giovanni Vigna, and Elaine Shi (eds.), *Proceedings of the ACM Conference on Computer and Communications (CCS)*, pp. 896–911. ACM, 2021.
- Min Chen, Zhikun Zhang, Tianhao Wang, Michael Backes, Mathias Humbert, and Yang Zhang. Graph unlearning. In *Proceedings of the ACM Conference on Computer and Communications Security (CCS)*, pp. 499–513, 2022.
- Min Chen, Weizhuo Gao, Gaoyang Liu, Kai Peng, and Chen Wang. Boundary unlearning: Rapid forgetting of deep networks via shifting the decision boundary. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7766–7775, 2023.
- Vikram S. Chundawat, Ayush K. Tarun, Murari Mandal, and Mohan S. Kankanhalli. Can bad teaching induce forgetting? unlearning in deep networks using an incompetent teacher. In Brian Williams, Yiling Chen, and Jennifer Neville (eds.), *Proceedings of the Thirty-Seventh Conference on Artificial Intelligence (AAAI)*, pp. 7210–7217, 2023.
- Benoît Colson, Patrice Marcotte, and Gilles Savard. An overview of bilevel optimization. *Annals of operations research*, 153:235–256, 2007.

- Vincent Conitzer and Tuomas Sandholm. Computing the optimal strategy to commit to. In Joan Feigenbaum, John C.-I. Chuang, and David M. Pennock (eds.), *Proceedings 7th ACM Conference on Electronic Commerce (EC)*, pp. 82–90. ACM, 2006.
- Steven Diamond and Stephen Boyd. Cvxpy: A python-embedded modeling language for convex optimization. *The Journal of Machine Learning Research*, 17(1):2909–2913, 2016.
- Asen L Dontchev, R Tyrrell Rockafellar, and R Tyrrell Rockafellar. *Implicit functions and solution mappings: A view from variational analysis*, volume 616. Springer, 2009.
- Aditya Golatkar, Alessandro Achille, and Stefano Soatto. Eternal sunshine of the spotless net: Selective forgetting in deep networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9304–9312, 2020a.
- Aditya Golatkar, Alessandro Achille, and Stefano Soatto. Eternal sunshine of the spotless net: Selective forgetting in deep networks. In *Proceedings of the Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 9301–9309, 2020b.
- Stephen Gould. Lecture notes on differentiable optimisation in deep learning. *CoRR*, 2023. URL <https://users.cecs.anu.edu.au/~sgould/papers/isaac22-lecture-notes.pdf>.
- Stephen Gould, Basura Fernando, Anoop Cherian, Peter Anderson, Rodrigo Santa Cruz, and Edison Guo. On differentiating parameterized argmin and argmax problems with application to bi-level optimization. *CoRR*, abs/1607.05447, 2016.
- Laura Graves, Vineel Nagisetty, and Vijay Ganesh. Amnesiac machine learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 11516–11524, 2021a.
- Laura Graves, Vineel Nagisetty, and Vijay Ganesh. Amnesiac machine learning. In *Proceedings of the 33rd Conference on Artificial Intelligence (AAAI)*, pp. 11516–11524. AAAI Press, 2021b.
- Chuan Guo, Tom Goldstein, Awni Y. Hannun, and Laurens van der Maaten. Certified data removal from machine learning models. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, volume 119, pp. 3832–3842. PMLR, 2020.
- Tomohiro Hayase, Suguru Yasutomi, and Takashi Katoh. Selective forgetting of deep networks at a finer level than samples. *arXiv preprint arXiv:2012.11849*, 2020.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, pp. 770–778, 2016.
- Zachary Izzo, Mary Anne Smart, Kamalika Chaudhuri, and James Zou. Approximate data deletion from machine learning models. In *International Conference on Artificial Intelligence and Statistics*, pp. 2008–2016. PMLR, 2021a.
- Zachary Izzo, Mary Anne Smart, Kamalika Chaudhuri, and James Zou. Approximate data deletion from machine learning models. In Arindam Banerjee and Kenji Fukumizu (eds.), *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 130, pp. 2008–2016, 2021b.
- Matthew Jagielski, Om Thakkar, Florian Tramèr, Daphne Ippolito, Katherine Lee, Nicholas Carlini, Eric Wallace, Shuang Song, Abhradeep Guha Thakurta, Nicolas Papernot, and Chiyuan Zhang. Measuring forgetting of memorized training examples. In *Proceedings the 11th International Conference on Learning Representations (ICLR)*, 2023.
- Jinghan Jia, Jiancheng Liu, Parikshit Ram, Yuguang Yao, Gaowen Liu, Yang Liu, Pranay Sharma, and Sijia Liu. Model sparsification can simplify machine unlearning. *CoRR*, abs/2304.04934, 2023. URL <https://doi.org/10.48550/arXiv.2304.04934>.
- Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning (ICML)*, volume 70, pp. 1885–1894, 2017a.

- Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In *International conference on machine learning*, pp. 1885–1894. PMLR, 2017b.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Meghdad Kurmanji, Peter Triantafillou, Jamie Hayes, and Eleni Triantafillou. Towards unbounded machine unlearning. *Advances in Neural Information Processing Systems*, 36, 2024.
- Frank J Massey Jr. The kolmogorov-smirnov test for goodness of fit. *Journal of the American statistical Association*, 46(253):68–78, 1951.
- Sasi Kumar Murakonda and Reza Shokri. Ml privacy meter: Aiding regulatory compliance by quantifying the privacy risks of machine learning. *CoRR*, abs/2007.09339, 2020.
- Milad Nasr, Reza Shokri, and Amir Houmansadr. Machine learning with membership privacy using adversarial regularization. In *Proceedings of the ACM Conference on Computer and Communications Security (CCS)*, pp. 634–646, 2018.
- Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011.
- 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. *CoRR*, abs/2209.02299, 2022.
- Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017.
- Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models. In *Proceedings of the IEEE Symposium on Security and Privacy (S&P)*, pp. 3–18. IEEE, 2017.
- Liwei Song and Prateek Mittal. Systematic evaluation of privacy risks of machine learning models. In *Proceedings of the 30th USENIX Security Symposium (USENIX)*, pp. 2615–2632, 2021.
- Anvith Thudi, Gabriel Deza, Varun Chandrasekaran, and Nicolas Papernot. Unrolling sgd: Understanding factors influencing machine unlearning. In *2022 IEEE 7th European Symposium on Security and Privacy (EuroS&P)*, pp. 303–319. IEEE, 2022a.
- Anvith Thudi, Hengrui Jia, Ilia Shumailov, and Nicolas Papernot. On the necessity of auditable algorithmic definitions for machine unlearning. In Kevin R. B. Butler and Kurt Thomas (eds.), *Proceedings of the Thirty-first USENIX Security Symposium (USENIX)*, pp. 4007–4022, 2022b.
- Liang Tong, Sixie Yu, Scott Alfeld, and Yevgeniy Vorobeychik. Adversarial regression with multiple learners. In *Proceedings of the 35th International Conference on Machine Learning (ICML)*, volume 80, pp. 4953–4961, 2018.
- Alexander Warnecke, Lukas Pirch, Christian Wressnegger, and Konrad Rieck. Machine unlearning of features and labels. *arXiv preprint arXiv:2108.11577*, 2021.
- Alexander Warnecke, Lukas Pirch, Christian Wressnegger, and Konrad Rieck. Machine unlearning of features and labels. In *Proceedings of the 30th Annual Network and Distributed System Security Symposium (NDSS)*, 2023.
- Jiayuan Ye, Aadyaa Maddi, Sasi Kumar Murakonda, Vincent Bindschaedler, and Reza Shokri. Enhanced membership inference attacks against machine learning models. In *Proceedings of the Conference on Computer and Communications Security (CCS)*, pp. 3093–3106. ACM, 2022.

## A APPENDIX

### A.1 NOTATION TABLE

Table 4: A summary of the notations used in the paper

Notation	Meaning
$\mathcal{D} = \{(x_i, y_i)\}$	A dataset
$(x_i, y_i)$	One data point where $x_i$ is the feature while $y_i$ is the label
$\mathcal{X}, \mathcal{Y}$	The feature space and the label space
$D_f, D_{te}, D_{val}, D_{tr}, D_r$	The forgetting, testing, validation, training, and retain set
$(x_j^f, y_j^f)$	One data point from the forget set $D_f$
$\mathcal{A}$	A machine learning algorithm
$\mathcal{U}$	A machine unlearning algorithm
$\theta_o$	The original model, i.e., $\mathcal{A}(D_{tr})$
$\theta_u$	The unlearned model, i.e., $\mathcal{U}(\theta_o)$
$\theta_r$	The retrained model, i.e., $\mathcal{A}(D_r)$
$\tilde{D}_{\theta_u}$	The auditing dataset for membership inference attack
$s_j^f$	The output of a forget instance in the auditing dataset $\tilde{D}_{\theta_u}$ from $\theta_u$
$s_j^{te}$	The output of a testing instance in the auditing dataset $\tilde{D}_{\theta_u}$ from $\theta_u$
$\tilde{D}_{\theta_u}^{tr}, \tilde{D}_{\theta_u}^{val}$	The training and validation split of $\tilde{D}_{\theta_u}$ .
$C_u$	The unlearner’s cost function
$\mathcal{B}_{\theta_u}$	The auditor’s best response given an unlearning model $\theta_u$
$U_a$	The utility function of the auditor
$\mathcal{H}_a$	The hypothesis class of the auditor
$\mathcal{H}_u$	The hypothesis class of the unlearner
$\alpha$	The trade-off factor as defined in the unlearner’s cost function equation 2

### A.2 RANDOM FORGETTING

#### A.2.1 SVHN DATASET

The results of SVHN dataset on random forgetting are given in Table 5.

Table 5: Experimental results (Mean<sub>std</sub>) on SVHN for random forgetting. The highlighted metrics are the closest to those of retraining, which is considered as the best performance compared with the other baselines.

SVHN	$Acc_r$	$Acc_{te}$	$Acc_f$	$ Acc_f - Acc_{te} $	MIA acc.	MIA AUC	MIA FI	KS Stat.	W. Dist.	RTE (min., ↓)
Retrain	0.9959 <sub>0.0002</sub>	0.9610 <sub>0.0010</sub>	0.9534 <sub>0.0024</sub>	0.0076	0.5248 <sub>0.0058</sub>	0.5422 <sub>0.0075</sub>	0.5149 <sub>0.0157</sub>	0.0306 <sub>0.0117</sub>	0.0686 <sub>0.0145</sub>	20.46
FT	0.9991 <sub>0.0001</sub>	0.7117 <sub>0.0021</sub>	0.9876 <sub>0.0070</sub>	0.2867	0.5372 <sub>0.0123</sub>	0.5592 <sub>0.0121</sub>	0.5523 <sub>0.0170</sub>	0.0613 <sub>0.0342</sub>	0.1743 <sub>0.0091</sub>	1.55
GA	0.9954 <sub>0.0001</sub>	0.9641 <sub>0.0002</sub>	0.9949 <sub>0.0006</sub>	0.0308	0.5191 <sub>0.0072</sub>	0.5411 <sub>0.0051</sub>	0.5500 <sub>0.0178</sub>	0.0867 <sub>0.0065</sub>	0.1473 <sub>0.0026</sub>	0.97
IU	0.9076 <sub>0.0707</sub>	0.8817 <sub>0.0658</sub>	0.9050 <sub>0.0713</sub>	0.0233	0.5373 <sub>0.0116</sub>	0.5580 <sub>0.0097</sub>	0.5469 <sub>0.0187</sub>	0.0473 <sub>0.0207</sub>	0.1407 <sub>0.0882</sub>	<b>0.41</b>
$\ell_1$ -sparse	0.9378 <sub>0.0615</sub>	0.9191 <sub>0.0540</sub>	0.9298 <sub>0.0620</sub>	0.0107	0.5457 <sub>0.0220</sub>	0.5665 <sub>0.0229</sub>	0.5347 <sub>0.0324</sub>	<b>0.0396</b> <sub>0.0112</sub>	0.1158 <sub>0.0954</sub>	1.86
RL	0.9949 <sub>0.0002</sub>	<b>0.9609</b> <sub>0.0006</sub>	0.9797 <sub>0.0018</sub>	0.0188	0.5211 <sub>0.0106</sub>	0.5411 <sub>0.0147</sub>	<b>0.5144</b> <sub>0.0225</sub>	0.1079 <sub>0.0175</sub>	<b>0.0642</b> <sub>0.0060</sub>	2.65
BE	0.9955 <sub>0.0001</sub>	0.9633 <sub>0.0002</sub>	0.9955 <sub>0.0006</sub>	0.0322	0.5209 <sub>0.0090</sub>	0.5441 <sub>0.0064</sub>	0.5553 <sub>0.0175</sub>	0.1016 <sub>0.0062</sub>	0.1528 <sub>0.0019</sub>	0.46
BS	<b>0.9956</b> <sub>0.0002</sub>	0.9641 <sub>0.0001</sub>	0.9952 <sub>0.0008</sub>	0.0311	0.5322 <sub>0.0060</sub>	0.5509 <sub>0.0034</sub>	0.5594 <sub>0.0176</sub>	0.0994 <sub>0.0074</sub>	0.1404 <sub>0.0033</sub>	0.81
SCRUB	0.9832 <sub>0.0010</sub>	0.9559 <sub>0.0014</sub>	0.9809 <sub>0.0020</sub>	0.0250	<b>0.5273</b> <sub>0.0031</sub>	<b>0.5431</b> <sub>0.0103</sub>	0.5296 <sub>0.0196</sub>	0.0492 <sub>0.0139</sub>	0.1032 <sub>0.0141</sub>	1.85
SG	0.9686 <sub>0.0017</sub>	0.9576 <sub>0.0033</sub>	<b>0.9560</b> <sub>0.0027</sub>	<b>0.0016</b>	0.5012 <sub>0.0052</sub>	0.5089 <sub>0.0272</sub>	0.3292 <sub>0.1798</sub>	0.0594 <sub>0.0233</sub>	0.0185 <sub>0.0041</sub>	3.16

### A.3 TINYIMAGENET DATASET

The results of TinyImageNet dataset on random forgetting are given in Table 6.

### A.4 CELEBA DATASET

The results of CelebA dataset on random forgetting are given in Table 7.

Table 6: Experimental results (Mean<sub>std</sub>) on TinyImageNet for random forgetting. The highlighted metrics are the closest to those of retraining, which is considered as the best performance compared with the other baselines.

TinyImageNet	$Acc_r$	$Acc_{te}$	$Acc_f$	$ Acc_f - Acc_{te} $	MIA acc.	MIA AUC	MIA FI	KS Stat.	W. Dist.	RTE (min., ↓)
Retrain	0.8377 <sub>0.0009</sub>	0.5967 <sub>0.0045</sub>	0.5057 <sub>0.0014</sub>	0.0910	0.5471 <sub>0.0028</sub>	0.5677 <sub>0.0029</sub>	0.4803 <sub>0.0037</sub>	0.1101 <sub>0.0021</sub>	0.4608 <sub>0.0124</sub>	237.12
FT	0.8242 <sub>0.0009</sub>	0.6095 <sub>0.0023</sub>	0.7033 <sub>0.0015</sub>	0.0938	0.5402 <sub>0.0019</sub>	0.5336 <sub>0.0013</sub>	0.6056 <sub>0.0024</sub>	0.0957 <sub>0.0030</sub>	0.5017 <sub>0.0071</sub>	65.07
GA	0.8132 <sub>0.0138</sub>	<b>0.5966</b> <sub>0.0061</sub>	0.8056 <sub>0.0170</sub>	0.2090	0.5966 <sub>0.0056</sub>	0.6032 <sub>0.0057</sub>	0.6619 <sub>0.0059</sub>	0.1968 <sub>0.0103</sub>	0.9195 <sub>0.0371</sub>	12.49
IU	0.8359 <sub>0.0010</sub>	0.6061 <sub>0.0001</sub>	0.8340 <sub>0.0033</sub>	0.2269	0.6051 <sub>0.0029</sub>	0.6150 <sub>0.0026</sub>	0.6708 <sub>0.0032</sub>	0.2170 <sub>0.0007</sub>	0.9684 <sub>0.0082</sub>	<b>6.73</b>
$\ell_1$ -sparse	0.7820 <sub>0.0015</sub>	0.6144 <sub>0.0012</sub>	0.6379 <sub>0.0045</sub>	0.0231	0.5039 <sub>0.0034</sub>	0.4906 <sub>0.0026</sub>	0.5674 <sub>0.0040</sub>	0.0500 <sub>0.0025</sub>	0.2217 <sub>0.0082</sub>	103.85
RL	0.7747 <sub>0.0006</sub>	0.6018 <sub>0.0020</sub>	0.5916 <sub>0.0030</sub>	0.0102	0.5280 <sub>0.0025</sub>	<b>0.5702</b> <sub>0.0018</sub>	0.4753 <sub>0.0047</sub>	0.1661 <sub>0.0013</sub>	0.3304 <sub>0.0034</sub>	60.79
BE	0.8054 <sub>0.0115</sub>	0.5561 <sub>0.0120</sub>	0.8038 <sub>0.0162</sub>	0.2477	0.6217 <sub>0.0026</sub>	0.6341 <sub>0.0033</sub>	0.6779 <sub>0.0005</sub>	0.2403 <sub>0.0047</sub>	1.0962 <sub>0.0040</sub>	31.37
BS	<b>0.8261</b> <sub>0.0008</sub>	0.5775 <sub>0.0001</sub>	0.8232 <sub>0.0020</sub>	0.2457	0.6234 <sub>0.0044</sub>	0.6333 <sub>0.0025</sub>	0.6818 <sub>0.0038</sub>	0.2415 <sub>0.0049</sub>	1.0623 <sub>0.0070</sub>	8.67
SG	0.8486 <sub>0.0007</sub>	0.5976 <sub>0.0005</sub>	<b>0.5560</b> <sub>0.0053</sub>	<b>0.0416</b>	<b>0.5212</b> <sub>0.0055</sub>	0.5341 <sub>0.0070</sub>	<b>0.5522</b> <sub>0.0198</sub>	<b>0.0893</b> <sub>0.0003</sub>	<b>0.3416</b> <sub>0.0046</sub>	13.36

Table 7: Experimental results (Mean<sub>std</sub>) on CelebA for random forgetting. The highlighted metrics are the closest to those of retraining, which is considered as the best performance compared with the other baselines.

CelebA	$Acc_r$	$Acc_{te}$	$Acc_f$	$ Acc_f - Acc_{te} $	MIA acc.	MIA AUC	MIA FI	KS Stat.	W. Dist.	RTE (min., ↓)
Retrain	0.9584 <sub>0.0114</sub>	0.9087 <sub>0.0110</sub>	0.9284 <sub>0.0048</sub>	0.0076	0.5123 <sub>0.0085</sub>	0.5103 <sub>0.0087</sub>	0.6385 <sub>0.0082</sub>	0.0285 <sub>0.0155</sub>	0.0686 <sub>0.0013</sub>	33.70
FT	0.9361 <sub>0.0010</sub>	0.9257 <sub>0.0021</sub>	0.9320 <sub>0.0028</sub>	0.0063	0.5038 <sub>0.0003</sub>	0.5070 <sub>0.0037</sub>	0.6180 <sub>0.0037</sub>	0.0133 <sub>0.0008</sub>	0.0158 <sub>0.0052</sub>	3.43
GA	0.9444 <sub>0.0003</sub>	0.9276 <sub>0.0001</sub>	0.9480 <sub>0.0012</sub>	0.0204	0.5215 <sub>0.0004</sub>	0.5214 <sub>0.0019</sub>	0.6325 <sub>0.0006</sub>	0.0236 <sub>0.0018</sub>	0.0491 <sub>0.0038</sub>	2.43
$\ell_1$ -sparse	0.7104 <sub>0.0883</sub>	0.7040 <sub>0.0917</sub>	0.7163 <sub>0.0921</sub>	0.0123	0.5038 <sub>0.0002</sub>	0.5074 <sub>0.0068</sub>	0.5850 <sub>0.0332</sub>	0.0251 <sub>0.0107</sub>	0.0277 <sub>0.0125</sub>	5.12
RL	0.9363 <sub>0.0007</sub>	0.9257 <sub>0.0005</sub>	0.9379 <sub>0.0003</sub>	0.0122	0.5036 <sub>0.0006</sub>	0.5081 <sub>0.0047</sub>	0.6238 <sub>0.0025</sub>	0.0188 <sub>0.0026</sub>	0.0343 <sub>0.0069</sub>	2.65
BE	0.9386 <sub>0.0027</sub>	<b>0.9210</b> <sub>0.0021</sub>	0.9406 <sub>0.0015</sub>	0.0196	0.5081 <sub>0.0005</sub>	0.5171 <sub>0.0038</sub>	0.6055 <sub>0.0009</sub>	<b>0.0289</b> <sub>0.0048</sub>	<b>0.0546</b> <sub>0.0049</sub>	3.59
BS	<b>0.9434</b> <sub>0.0005</sub>	0.9270 <sub>0.0002</sub>	0.9465 <sub>0.0043</sub>	0.0195	0.5096 <sub>0.0012</sub>	0.5139 <sub>0.0005</sub>	0.6287 <sub>0.0018</sub>	0.0270 <sub>0.0008</sub>	0.0483 <sub>0.0065</sub>	<b>1.67</b>
SG	0.9348 <sub>0.0011</sub>	0.9222 <sub>0.0001</sub>	<b>0.9288</b> <sub>0.0010</sub>	<b>0.0066</b>	<b>0.5159</b> <sub>0.0002</sub>	<b>0.5103</b> <sub>0.0007</sub>	<b>0.6358</b> <sub>0.0009</sub>	0.0274 <sub>0.0003</sub>	0.0108 <sub>0.0006</sub>	9.74

## A.5 ViT RESULTS

The comparison of SG on SVHN for random forgetting with and without attacker is illustrated in Figure 3.

### A.5.1 LOSS DISTRIBUTIONS

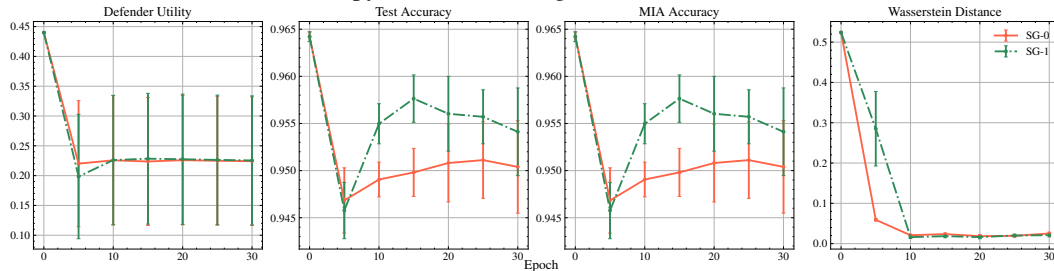
A visualization of the cross-entropy losses of the forget and test instances is in Figure 4.

## A.6 BASELINE METHODS

**Retrain:** The first baseline is retraining, where the unlearned model is obtained by training on the retain set from scratch. We aim to develop unlearning algorithms so that the metrics they produce are as closely aligned with those of the retraining as possible.

**Fine-Tuning (FT):** As the second baseline, FT continues to train the original model on the retain set for a few epochs. This a standard baseline used in various prior research (Graves et al., 2021b; Warnecke et al., 2023).

Figure 3: An ablation study to understand the impact of adversarial modeling on the process of unlearning;  $\alpha = 1$  and  $\alpha = 0$  corresponds to the cases with and without adversarial modeling, respectively. The results are the averages over 10 experiments with different seeds, and 95% confidence intervals are displayed. **From the left to the right:** 1) the defender’s utility, evaluated as the test accuracy  $Acc_{te}$  minus the MIA accuracy; 2) test accuracy; 3) MIA accuracy; 4) Wasserstein distance between the cross-entropy losses of the forget and test instances.



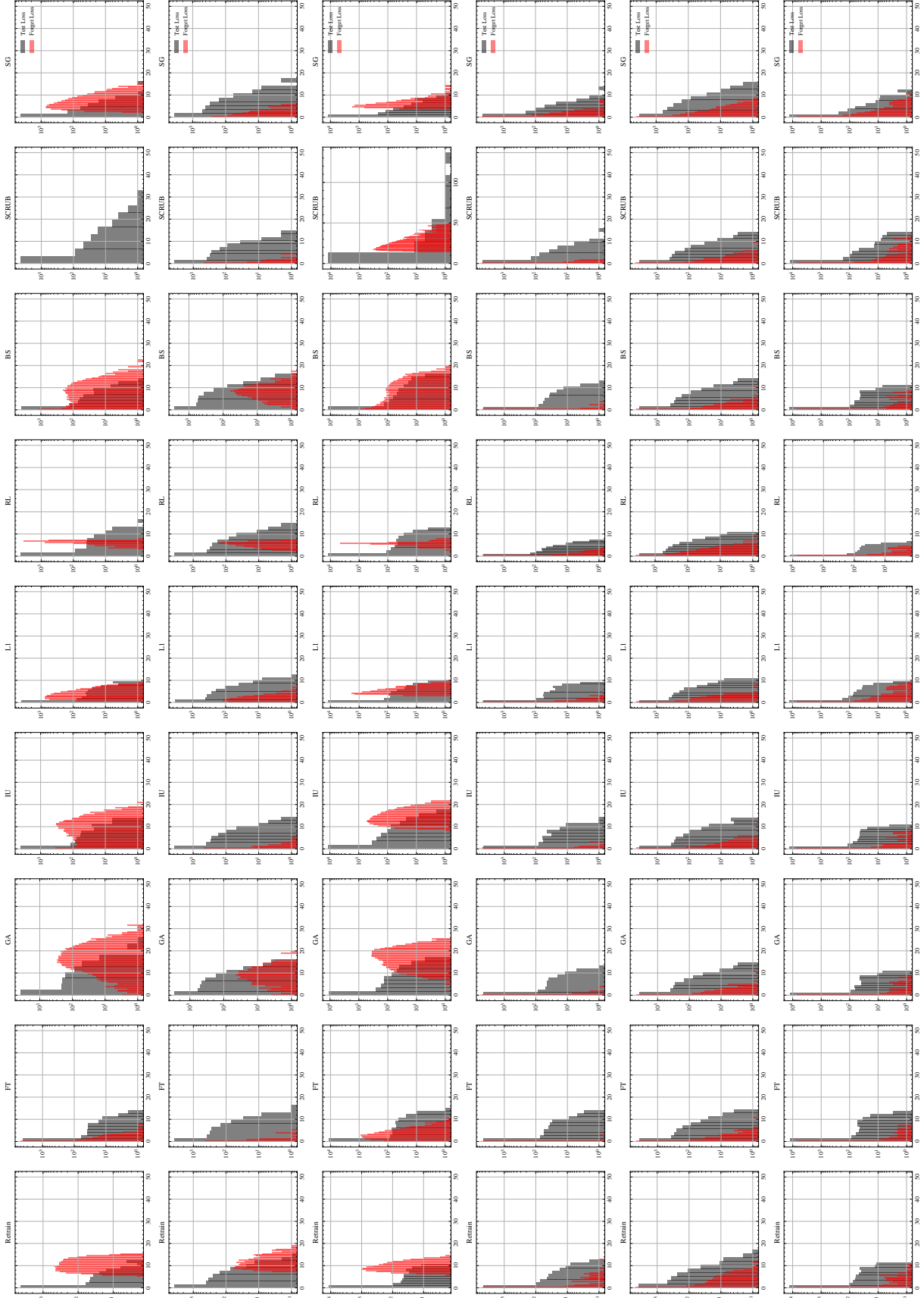


Figure 4: The distributions of the cross-entropy losses for the forget and test instances from the unlearned models. The y-axis is in log scale for better visualization. From the first to the last figure, they are random forgetting on CIFAR-10, CIFAR-100, SVHN and class-wise forgetting on CIFAR-10, CIFAR-100, SVHN.



Table 8: Experimental results (Mean<sub>std</sub>) on CIFAR-10 for random forgetting using ViT. The high-lighted metrics are the closest to those of retraining, which is considered as the best performance compared with the other baselines.

CIFAR-10	$Acc_r$	$Acc_{te}$	$Acc_f$	$ Acc_f - Acc_{te} $	MIA acc.	MIA AUC	MIA FI	KS Stat.	W. Dist.	RTE (min., ↓)
Retrain	0.8384 <sub>0.0020</sub>	0.7427 <sub>0.0017</sub>	0.7483 <sub>0.0033</sub>	0.0056	0.5000 <sub>0.0007</sub>	0.5000 <sub>0.0059</sub>	0.5794 <sub>0.0163</sub>	0.0190 <sub>0.0105</sub>	0.0304 <sub>0.0073</sub>	206.50s
FT	0.8630 <sub>0.0017</sub>	0.7599 <sub>0.0033</sub>	0.8221 <sub>0.0083</sub>	0.0622	0.5286 <sub>0.0001</sub>	0.5352 <sub>0.0005</sub>	0.6188 <sub>0.0019</sub>	0.0675 <sub>0.0078</sub>	0.2413 <sub>0.0142</sub>	18.85
GA	0.8473 <sub>0.0017</sub>	0.7594 <sub>0.0105</sub>	0.8461 <sub>0.0018</sub>	0.0867	0.5418 <sub>0.0020</sub>	0.5512 <sub>0.0050</sub>	0.6302 <sub>0.0002</sub>	0.0898 <sub>0.0074</sub>	0.3055 <sub>0.0313</sub>	<b>3.01</b>
$\ell_1$ -sparse	0.8472 <sub>0.0015</sub>	0.7591 <sub>0.0103</sub>	0.8457 <sub>0.0007</sub>	0.0866	0.5423 <sub>0.0016</sub>	0.5512 <sub>0.0050</sub>	0.6305 <sub>0.0002</sub>	0.0890 <sub>0.0076</sub>	0.3048 <sub>0.0313</sub>	6.32
RL	<b>0.8415</b> <sub>0.0027</sub>	0.7592 <sub>0.0085</sub>	<b>0.8124</b> <sub>0.0034</sub>	<b>0.0532</b>	0.5157 <sub>0.0008</sub>	0.5114 <sub>0.0056</sub>	0.5862 <sub>0.0017</sub>	0.0598 <sub>0.0099</sub>	0.1473 <sub>0.0158</sub>	8.09
SG	0.8515 <sub>0.0045</sub>	<b>0.7476</b> <sub>0.0203</sub>	0.8322 <sub>0.0091</sub>	0.0846	<b>0.5019</b> <sub>0.0074</sub>	<b>0.5100</b> <sub>0.0013</sub>	<b>0.5814</b> <sub>0.0010</sub>	<b>0.0374</b> <sub>0.0004</sub>	<b>0.1041</b> <sub>0.0034</sub>	11.42

Table 9: The hyper-parameter for the baseline method and SG used in this paper.

Parameters	Retrain	FT	GA	IU	$\ell_1$ -sparse	RL	BE	BS	SCRUB	SG
Learning rate	1e-2	5e-2	1e-3	×	1e-2	1e-2	1e-5	1e-5	5e-4	1e-2
Num. of epoch	160	30	5	×	10	10	10	10	10	30
$\gamma$	×	×	×	×	5e-4	×	×	×	×	×
$\alpha$	×	×	×	10	×	×	×	×	×	×
$T$	×	×	×	×	×	×	×	×	4	×
Decay epochs	×	×	×	×	×	×	×	×	[3, 5, 9]	×
$\beta$	×	×	×	×	×	×	×	×	0.1	×
Attacker $\alpha$	×	×	×	×	×	×	×	×	×	1.0

**Gradient Ascent (GA):** This baseline takes the original model as the starting point and runs a few epochs of gradient ascent on the forget set  $D_f$ . The intuition is to disrupt the model’s generalizability on  $D_f$  (Graves et al., 2021b). Another name of GA is NegGrad (Kurmanji et al., 2024).

**Influence Unlearning (IU):** This baseline uses Influence Function to estimate the updates required for a model’s weights as a result of removing the forget set from the training data (Izzo et al., 2021b; Koh & Liang, 2017a).

**$\ell_1$ -sparse:** This baseline integrates an  $\ell_1$  norm-based sparse penalty into machine unlearning loss Jia et al. (2023).

**Random Label (RL):** This baseline trains the original model on the retain set and the forgetting set  $D_f$  whose labels are random to make the model unlearn  $D_f$  while keep the model capability as much as possible.

**Boundary Expansion:** This baseline proposes a neighbor searching method to identify the nearest but incorrect class labels to guide the way of boundary shifting.

**Boundary Shrink:** This baseline artificially assigns forgetting samples to an extra shadow class of the original model Chen et al. (2023).

**SCRUB:** This baseline achieve MU by using a teacher model and student model Kurmanji et al. (2024).

## A.7 EXPERIMENT DETAILS

The hyperparameters used for SG and the baselines are in Table 9. The losses for the retraining baseline across the epochs are displayed in Figure 5. We run all the experiments using PyTorch 1.12 on NVIDIA A5000 GPUs and AMD EPYC 7513 32-Core Processor.

## A.8 CLASS-WISE FORGETTING

The results of SVHN dataset on classwise forgetting are given in Table 10.

## A.9 AN EXAMPLE OF THE CONDITION IN EQUATION 5

In this section, we provide a concrete example of the KKT conditions for linear support vector machines (SVM). As described in Section 1, the KKT conditions are key to relating the attacker’s model parameters, denoted as  $\theta_a$ , with the auditing set  $\tilde{D}_{\theta_u}$ , which allows us to derive the gradient  $\partial\theta_a/\partial\tilde{D}_{\theta_u}$ . The conditions  $f$  can be similarly derived for any model where the learning problem is convex. To simplify the notations, we use  $\{(x_i, y_i)\}_{i=1}^q$  to represent  $\tilde{D}_{\theta_u}$ . A standard formulation of

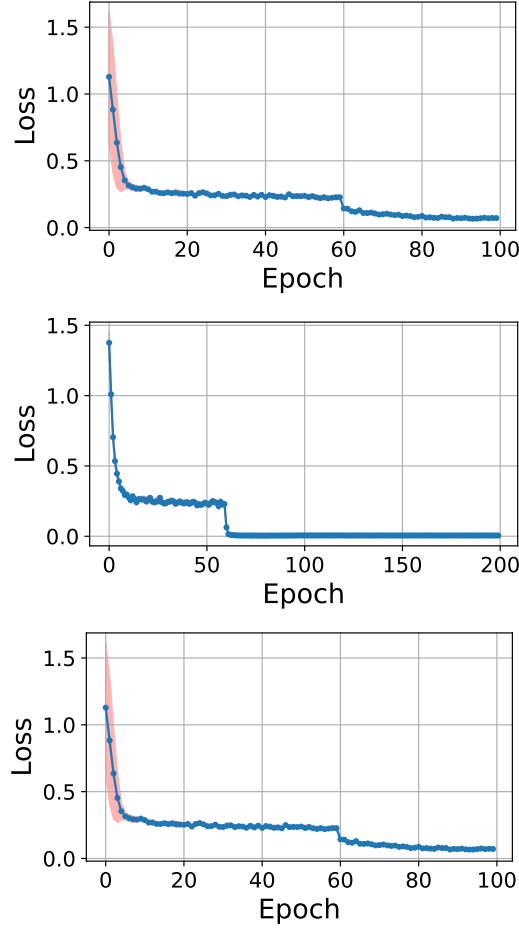


Figure 5: The training loss for the retrain baseline. For CIFAR10 and CIFAR100, the learning rate is multiplied by 0.1 when epoch is at 60, 120, 160; for SVHN, the same multiplication is done at epoch 60, 120. **Top to bottom:** CIFAR10, CIFAR100, SVHN.

Table 10: Experimental results (Mean<sub>std</sub>) on SVHN for classwise forgetting. The highlighted metrics are the closest to those of retraining, which is considered as the best performance compared with the other baselines.

SVHN	$Acc_r$	$Acc_{te}$	$Acc_f$	$ Acc_f - Acc_{te} $	MIA acc.	MIA AUC	MIA F1	KS Stat.	W. Dist.	RTE (min., ↓)
Retrain	0.9963 <sub>0.0001</sub>	0.9639 <sub>0.0007</sub>	0.0000 <sub>0.0000</sub>	0.9639	0.9950 <sub>0.0005</sub>	0.9986 <sub>0.0004</sub>	0.9951 <sub>0.0005</sub>	0.9909 <sub>0.0003</sub>	8.6924 <sub>0.0750</sub>	20.46
FT	<b>0.9978</b> <sub>0.0002</sub>	<b>0.9622</b> <sub>0.0020</sub>	0.0995 <sub>0.0179</sub>	0.8627	<b>0.9945</b> <sub>0.0005</sub>	<b>0.9985</b> <sub>0.0006</sub>	<b>0.9946</b> <sub>0.0007</sub>	0.9533 <sub>0.0038</sub>	2.3220 <sub>0.0216</sub>	3.17
GA	0.9444 <sub>0.0055</sub>	0.9144 <sub>0.0047</sub>	<b>0.0000</b> <sub>0.0000</sub>	0.9144	0.9969 <sub>0.0004</sub>	0.9998 <sub>0.0000</sub>	0.9970 <sub>0.0002</sub>	0.9849 <sub>0.0012</sub>	16.4834 <sub>0.2720</sub>	0.98
IU	0.8044 <sub>0.1177</sub>	0.8061 <sub>0.0978</sub>	<b>0.0000</b> <sub>0.0000</sub>	0.8061	0.9999 <sub>0.0003</sub>	1.0000 <sub>0.0000</sub>	0.9998 <sub>0.0003</sub>	<b>0.9936</b> <sub>0.0056</sub>	15.0697 <sub>1.7117</sub>	0.41
$\ell_1$ -sparse	0.9799 <sub>0.0004</sub>	0.9580 <sub>0.0017</sub>	<b>0.0000</b> <sub>0.0000</sub>	0.9580	0.9921 <sub>0.0013</sub>	0.9966 <sub>0.0003</sub>	0.9921 <sub>0.0012</sub>	0.9818 <sub>0.0030</sub>	4.5139 <sub>0.2512</sub>	3.73
RL	0.9959 <sub>0.0001</sub>	0.9612 <sub>0.0013</sub>	<b>0.0000</b> <sub>0.0000</sub>	0.9612	0.9912 <sub>0.0013</sub>	0.9971 <sub>0.0016</sub>	0.9913 <sub>0.0012</sub>	0.9813 <sub>0.0016</sub>	<b>5.5978</b> <sub>0.0357</sub>	2.60
BE	0.9880 <sub>0.0008</sub>	0.9546 <sub>0.0012</sub>	0.2812 <sub>0.0061</sub>	0.6734	0.9976 <sub>0.0005</sub>	0.9995 <sub>0.0002</sub>	0.9976 <sub>0.0006</sub>	0.9106 <sub>0.0050</sub>	4.3816 <sub>0.0659</sub>	0.46
BS	0.9864 <sub>0.0010</sub>	0.9537 <sub>0.0010</sub>	0.3109 <sub>0.0052</sub>	0.6428	0.9975 <sub>0.0003</sub>	0.9995 <sub>0.0002</sub>	0.9976 <sub>0.0003</sub>	0.9072 <sub>0.0031</sub>	4.4290 <sub>0.1169</sub>	0.82
SCRUB	0.9916 <sub>0.0007</sub>	0.9616 <sub>0.0014</sub>	<b>0.0000</b> <sub>0.0000</sub>	<b>0.9616</b>	0.9999 <sub>0.0001</sub>	1.0000 <sub>0.0001</sub>	0.9999 <sub>0.0001</sub>	0.9989 <sub>0.0008</sub>	24.4590 <sub>2.2852</sub>	3.91
SG	0.9716 <sub>0.0007</sub>	0.9601 <sub>0.0014</sub>	<b>0.0000</b> <sub>0.0000</sub>	0.9601	0.9928 <sub>0.0001</sub>	0.9954 <sub>0.0001</sub>	0.9929 <sub>0.0001</sub>	0.9907 <sub>0.0008</sub>	5.0148 <sub>2.2852</sub>	5.92

the linear SVM is as follows

$$\begin{aligned}
 \min_{\theta_a, b} \quad & \frac{1}{2} \|\theta_a\|^2 \\
 s.t. \quad & y_i \cdot (\theta_a^\top x_i + b) \geq 1, \forall i,
 \end{aligned} \tag{7}$$

where  $b$  is the bias term. The standard form is typically formulated as a minimization problem, so the attacker is to maximize  $V = -\frac{1}{2} \|\theta_a\|^2$ . Eq. equation 7 is a convex program, and the optimal solution (i.e.,  $\theta_a^*$  and  $b^*$ ) is characterized by the KKT conditions. The Lagrangian of the above is as follows

where  $\alpha_i \geq 0$  are the Lagrangian multipliers:

$$L(\theta_a, b, \alpha_i) = \frac{1}{2} \|\theta_a\|^2 - \sum_{i=1}^q \alpha_i (y_i \cdot (\theta_a^\top x_i + b) - 1). \quad (8)$$

Following standard procedures (Boyd & Vandenberghe, 2004), the KKT conditions are as follows

$$f(\tilde{D}_{\theta_u}, \theta_a) = \begin{cases} \theta_a - \sum_{i=1}^q \alpha_i y_i x_i = 0 \\ - \sum_{i=1}^q \alpha_i y_i = 0 \\ y_i \cdot (\theta_a^\top x_i + b) \geq 1 \\ \alpha_i \geq 0, \forall i \\ \alpha_i (y_i (\theta_a^\top x_i + b) - 1) = 0, \forall i \end{cases}, \quad (9)$$

which implicitly define a function between  $\theta_a$  and the data  $\tilde{D}_{\theta_u} = \{(x_i, y_i)\}_{i=1}^q$ . In practice, we describe the optimization problem equation 7 using `cvxpy` (Diamond & Boyd, 2016). Then, we employ an off-the-shelf package called `cvxpylayers` (Agrawal et al., 2019) to automatically derive the KKT conditions and compute the gradient  $\partial \theta_a / \partial \tilde{D}_{\theta_u}$ .