

# LABEL SMOOTHING IMPROVES MACHINE UNLEARNING

Zonglin Di<sup>1</sup>, Zhaowei Zhu<sup>2</sup>, Jinghan Jia<sup>3\*</sup>, Jiancheng Liu<sup>3\*</sup>, Zafar Takhirov<sup>4</sup>, Bo Jiang<sup>4</sup>, Yuanshun Yao<sup>5</sup>, Sijia Liu<sup>3</sup>, Yang Liu<sup>1†</sup>

<sup>1</sup>University of California, Santa Cruz, <sup>2</sup>Docta.ai, <sup>3</sup>Michigan State University,

<sup>4</sup>Independent Contributor, <sup>5</sup>Meta SuperIntelligence Lab

<sup>1</sup>{zdi, yangliu}@ucsc.edu, <sup>2</sup>zzw@docta.ai,

<sup>3</sup>{jiajingh, liujia45, liusiji5}@msu.edu,

<sup>4</sup>{z.takhirov, bjiang518}@gmail.com, <sup>5</sup>kevinyaowork@gmail.com

## ABSTRACT

The objective of machine unlearning (MU) is to eliminate previously learned data from a model. However, it can be challenging to strike a balance between computation cost and performance when using existing MU techniques. Taking inspiration from the influence of label smoothing on model confidence and differential privacy, we propose a simple gradient-based MU approach that uses an inverse process of label smoothing. This work introduces UGradSL, a simple, plug-and-play MU approach that uses smoothed labels. We provide theoretical analyses demonstrating why properly introducing label smoothing improves MU performance. We conducted extensive experiments on several datasets of various sizes and different modalities, demonstrating the effectiveness and robustness of our proposed method. UGradSL is also closely connected to improving the local differential privacy. The consistent improvement in MU performance is only at a marginal cost of additional computations. For instance, UGradSL consistently outperforms the gradient-ascent MU baseline on different unlearning tasks without sacrificing unlearning efficiency. A self-adaptive UGradSL is also given for simple parameter selection. The code is available at <https://github.com/UCSC-REAL/Label-Smoothing-Unlearn>.

## 1 INTRODUCTION

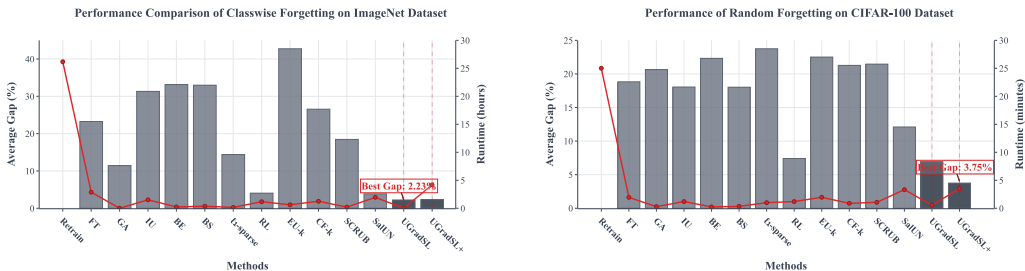
Building a reliable ML model has become an important topic in this community. Machine unlearning (MU) is a task requiring removing the learned data points from the model. The concept and the technology of MU enable researchers to delete sensitive or improper data in the training set to improve fairness, robustness, and privacy and get a better ML model for real-world use (Chen et al., 2021b; Sekhari et al., 2021). Retraining from scratch (Retrain) is a straightforward method when we want to remove the data from the model; yet it incurs prohibitive computation costs for large models due to computing resource constraints. Therefore, an efficient and effective MU method is desired.

The most straightforward MU approach should be a retraining-based method (Bourtoule et al., 2021), meaning that we retrain the model from scratch without using the data to be forgotten. The method can guarantee privacy protection but the computational cost is intensive. Most existing works (Koh & Liang, 2017; Golatkar et al., 2020; Warnecke et al., 2021; Graves et al., 2021; Thudi et al., 2021; Izzo et al., 2021; Becker & Liebig, 2022; Jia et al., 2023) focus on *approximate MU* to achieve a balance between unlearning efficacy and computational complexity, making them more suitable for real-world applications, meaning that they make the model unlearn the forgetting dataset without retraining the model.

We desire an approach that enjoys both high performance and high efficiency. Since MU can be viewed as the inverse process of ML, we are motivated to think it would be a natural and efficient

\*Equal Contribution.

†Corresponding author: Yang Liu <yangliu@ucsc.edu>



(a) Performance of classwise forgetting on ImageNet. (b) Performance of random forgetting on CIFAR-100.

Figure 1: The performance comparison of our proposed methods and baseline methods using average gap and runtime (RTE), where lower values indicate better performance. Bars represent average gap while red dotted lines show RTE. Since retraining does not have gap by definition, only RTE is reported for this baseline and the bar is empty. For classwise forgetting on ImageNet, UGradSL achieves the lowest average gap (2.23%) with acceptable RTE increase. For random forgetting on CIFAR-100, UGradSL+ attains the best average gap (3.75%), while UGradSL demonstrates an optimal gap-runtime trade-off.

way to develop an unlearning process that imitates the reverse of gradient descent. Indeed, gradient ascent (GA) (Thudi et al., 2021) is one of the MU methods but unfortunately, it does not fully achieve the potential of this idea. One of the primary reasons is that once the model completes training, the gradient of well-memorized data that was learned during the process diminishes (close to 0 loss) and therefore the effect of GA is rather limited.

Our approach is inspired by the celebrated idea of label smoothing (Szegedy et al., 2016). In the forward problem (gradient descent), the smoothed label proves to be able to improve the model’s generalization power. In our setting, we treat the smoothed label term as the regularization in the loss function, making the unlearning more controllable. Specifically, we show that GA with a “negative” label smoothing process (which effectively results in a standard label smoothing term in a descending fashion) can quickly improve the model’s deniability in the forgetting dataset, making the model behave close to the retrained model, which is exactly the goal of MU. We name our approach *UGradSL*, Unlearning using **Gradient-based Smoothed Labels**.

Our approach is a plug-and-play method that can improve the gradient-based MU performance consistently and does not hurt the performance of the remaining dataset and the testing dataset in a gradient-mixed way. At the same time, we provide a theoretical analysis of the benefits of our approach for the MU task. The core contributions of this paper are summarized as follows:

- We propose a lightweight tool to improve MU by joining the label smoothing and gradient ascent.
- We theoretically analyze the role of gradient ascent in MU and how negative label smoothing is able to boost MU performance.
- Extensive experiments in six datasets in different modalities and several unlearning paradigms regarding different MU metrics show the robustness and generalization of our method.
- We investigate the relationship between label smoothing and label differential privacy (LDP), showing that label smoothing can aid LDP.

## 2 RELATED WORK

**Machine Unlearning (MU)** was developed to address information leakage concerns related to private data after the completion of model training (Cao & Yang, 2015; Bourtole et al., 2021; Nguyen et al., 2022), gained prominence with the advent of privacy-focused legislation (Hoofnagle et al., 2019; Pardo, 2018). One direct unlearning method involves retraining the model from scratch after removing the forgetting data from the original training set. It is computationally inefficient, prompting researchers to focus on developing approximate but much faster unlearning techniques (Becker & Liebig, 2022; Golatkar et al., 2020; Warnecke et al., 2021; Graves et al., 2021; Thudi et al., 2021; Izzo et al., 2021; Jia et al., 2023). Beyond unlearning methods, other research efforts aim to create probabilistic unlearning concepts (Ginart et al., 2019; Guo et al., 2019; Neel et al.,

2021; Ullah et al., 2021; Sekhari et al., 2021) and facilitate unlearning with provable error guarantees, particularly in the context of differential privacy (DP) (Dwork et al., 2006; Ji et al., 2014; Hall et al., 2012). However, it typically necessitates stringent model and algorithmic assumptions, potentially compromising effectiveness against practical adversaries, such as membership inference attacks (Graves et al., 2021; Thudi et al., 2021). Additionally, the interest in MU has expanded to encompass various learning tasks and paradigms (Wang et al., 2022b; Liu et al., 2022b; Chen et al., 2022; Chien et al., 2022; Marchant et al., 2022; Di et al., 2022). These applications demonstrate the growing importance of MU techniques in safeguarding privacy. The rest of the related work about influence function label smoothing and differential privacy are given in Appendix B.

### 3 LABEL SMOOTHING ENABLES FAST AND EFFECTIVE UNLEARNING

This section sets up the analysis and shows that properly performing label smoothing enables fast and effective unlearning. The key ingredients of our approach are gradient ascent (GA) and label smoothing (LS). We start with understanding how GA helps with unlearning and then move on to show the power of LS. At the end of the section, we formally present our algorithm.

#### 3.1 PRELIMINARY

**Machine Unlearning** Consider a  $K$ -class classification problem on the training data distribution  $\mathcal{D}_{tr} = (\mathcal{X} \times \mathcal{Y})$ , where  $\mathcal{X}$  and  $\mathcal{Y}$  are the feature and label space, respectively. Due to some privacy regulations, there exists a forgetting data distribution  $\mathcal{D}_f$  that the model needs to unlearn. We denote by  $\theta_{tr}$  the original model trained on  $\mathcal{D}_{tr}$  and  $\theta_u$  the model without the influence of  $\mathcal{D}_f$ . The goal of machine unlearning (MU) is how to generate  $\theta_u$  from  $\theta_{tr}$ .

**Label Smoothing** In a  $K$ -class classification task, let  $\mathbf{y}_i$  denote the one-hot encoded vector form of  $y_i \in \mathcal{Y}$ . Similar to (Wei et al., 2021), we unify positive label smoothing (PLS) and negative label smoothing (NLS) into generalized label smoothing (GLS). The random variable of smoothed label  $\mathbf{y}_i^{\text{GLS}, \alpha}$  with smooth rate  $\alpha \in (-\infty, 1]$  is  $\mathbf{y}_i^{\text{GLS}, \alpha} = (1 - \alpha) \cdot \mathbf{y}_i + \frac{\alpha}{K} \cdot \mathbf{1} = [\frac{\alpha}{K}, \dots, \frac{\alpha}{K}, (1 + \frac{1-K}{K}\alpha), \frac{\alpha}{K}, \dots, \frac{\alpha}{K}]$ , where  $(1 + \frac{1-K}{K}\alpha)$  is the  $y_i$ -th element in the encoded label vector. When  $\alpha < 0$ , GLS becomes NLS.

#### 3.2 GRADIENT ASCENT CAN HELP GRADIENT-BASED MACHINE UNLEARNING

We discuss three sets of model parameters in the MU problem: 1)  $\theta_{tr}^*$ , the optimal parameters trained from  $D_{tr} \sim \mathcal{D}_{tr}$ , 2)  $\theta_r^*$ , the optimal parameters trained from  $D_r \sim \mathcal{D}_r$ , such that  $D_r = D_{tr} \setminus D_f$  and 3)  $\theta_f^*$ , the optimal parameters unlearned using gradient ascent (GA) on  $D_f \sim \mathcal{D}_f$ . Note  $\theta_r^*$  can be viewed as the *exact* MU model. The definitions of  $\theta_{tr}^*$  and  $\theta_r^*$  follow the standard empirical risk minimization as

$$\theta^* = \arg \min_{\theta} \frac{1}{n} \sum_{z \in D} \ell(h_{\theta}, z). \quad (1)$$

and by using the influence function,  $\theta_f^*$  is

$$\theta_f^* = \arg \min_{\theta} \{R_{tr}(\theta) + \varepsilon \sum_{z^f \in D_f} \ell(h_{\theta}, z^f)\}$$

where  $R_{tr}(\theta) = \sum_{z^{tr} \in D_{tr}} \ell(h_{\theta}, z^{tr})$  and  $R_f(\theta) = \sum_{z^f \in D_f} \ell(h_{\theta}, z^f)$  are the empirical risk on  $D_{tr}$  and  $D_f$ , respectively. We use notations  $\ell(h_{\theta}, z)$  to specify the loss of an example  $z = (x, y)$  in the dataset.  $h_{\theta}$  is a function  $h$  parameterized by  $\theta$ .  $\varepsilon$  is the weight of  $D_f$  compared with  $D_{tr}$ . The optimal parameter can be found when the gradient is 0:

$$\nabla_{\theta} R_{tr}(\theta_f^*) + \varepsilon \sum_{z^f \in D_f} \nabla_{\theta} \ell(h_{\theta_f^*}, z^f) = 0. \quad (2)$$

Expanding Eq. (2) at  $\theta = \theta_{tr}^*$  using the Taylor series, we have

$$\theta_f^* - \theta_{tr}^* \approx - \left[ \sum_{z^{tr} \in D_{tr}} \nabla_{\theta}^2 \ell(h_{\theta_{tr}^*}, z^{tr}) + \varepsilon \sum_{z^f \in D_f} \nabla_{\theta}^2 \ell(h_{\theta_{tr}^*}, z^f) \right]^{-1} \left( \varepsilon \sum_{z^f \in D_f} \nabla_{\theta} \ell(h_{\theta_{tr}^*}, z^f) \right). \quad (3)$$

Here, we ignore the Lagrange Remainder. Similarly, we can expand  $\nabla_{\theta} R_{tr}(\theta_{tr}^*)$  at  $\theta = \theta_r^*$  and derive  $\theta_r^* - \theta_{tr}^*$  as

$$\theta_r^* - \theta_{tr}^* \approx \left[ \sum_{z^{tr} \in D_{tr}} \nabla_{\theta}^2 \ell(h_{\theta_r^*}, z^{tr}) \right]^{-1} \left( \sum_{z^{tr} \in D_{tr}} \nabla_{\theta} \ell(h_{\theta_r^*}, z^{tr}) \right). \quad (4)$$

We ignore the average operation in the original definition of the influence function for computation convenience because the size of  $D_{tr}$  or  $D_f$  is fixed. For GA, let  $\varepsilon = -1$  in Eq. (3) and we have

$$\theta_r^* - \theta_f^* = \theta_r^* - \theta_{tr}^* - (\theta_f^* - \theta_{tr}^*) = \Delta\theta_r - \Delta\theta_f, \quad (5)$$

where  $(-\Delta\theta_r)$  represents the learning gap from  $\theta_r^*$  to  $\theta_{tr}^*$  while vector  $\Delta\theta_f$  represents how much the model unlearns (backtracked progress) between  $\theta_f^*$  and  $\theta_{tr}^*$ . The details of  $\Delta\theta_r$  and  $\Delta\theta_f$  are given in Eq. (17) in Appendix. Ideally, when  $\Delta\theta_r$  and  $\Delta\theta_f$  are exactly the same vectors, GA can lead the model to the optimal retrained model since we have  $\theta_r^* = \theta_f^*$ . However, this condition is hard to satisfy in practice. Thus, GA cannot always help MU. We summarize it in Theorem 1. The proof and the error analysis are given in Appendix C.1 and C.2.

**Theorem 1.** *Given the approximation in Eq. (5), GA achieves exact MU if and only if*

$$\sum_{z^f \in D_f} \nabla_{\theta} \ell(h_{\theta_r^*}, z^f) \approx -\mathbf{H}(\theta_r^*, \theta_{tr}^*) \cdot \sum_{z^f \in D_f} \nabla_{\theta} \ell(h_{\theta_{tr}^*}, z^f),$$

$\mathbf{H}(\theta_r^*, \theta_{tr}^*) = \left[ \sum_{z^{tr} \in D_{tr}} \nabla_{\theta}^2 \ell(h_{\theta_r^*}, z^{tr}) \right] \left[ \sum_{z^r \in D_r} \nabla_{\theta}^2 \ell(h_{\theta_{tr}^*}, z^r) \right]^{-1}$ . Otherwise, there exist  $\theta_r^*, \theta_{tr}^*$  such that GA cannot help MU, i.e.,  $\|\theta_r^* - \theta_f^*\| > \|\theta_r^* - \theta_{tr}^*\|$ .

### 3.3 LABEL SMOOTHING IMPROVES MU

Practically, we cannot guarantee that GA always helps MU as shown in Theorem 1. To alleviate the possible undesired effect of GA, we propose to use label smoothing as a plug-in module. Consider the cross-entropy loss as an example. For GLS, the loss is calculated as

$$\ell(h_{\theta}, z^{\text{GLS}, \alpha}) = \left(1 + \frac{1-K}{K}\alpha\right) \cdot \ell(h_{\theta}, (x, y)) + \frac{\alpha}{K} \sum_{y' \in \mathcal{Y} \setminus y} \ell(h_{\theta}, (x, y')), \quad (6)$$

where  $\ell(h_{\theta}, (x, y)) := \ell(h_{\theta}, z)$  and  $\ell(h_{\theta}, (x, y'))$  to denote the loss of an example when its label is replaced with  $y'$ . Intuitively, Term  $\sum_{y' \in \mathcal{Y} \setminus y} \ell(h_{\theta}, (x, y'))$  in Eq. (6) leads to a state where the model makes *wrong predictions on data in the forgetting dataset* with equally low confidence (Wei et al., 2021; Lukasik et al., 2020).

With smoothed label given in Eq. (6), we show that there exists a vector  $\Delta\theta_n$  such that Eq. (5) can be written as

$$\theta_r^* - \theta_{f, \text{LS}}^* \approx \Delta\theta_r - \Delta\theta_f + \frac{1-K}{K}\alpha \cdot (\Delta\theta_n - \Delta\theta_f), \quad (7)$$

We leave the detailed form of  $\Delta\theta_n$  to Eq. (34). But intuitively,  $\Delta\theta_n$  captures the gradient influence of the smoothed non-target label on the weight. We show the effect of NLS ( $\alpha < 0$ ) in Theorem 2 below and its proof is given in Appendix C.3.

**Theorem 2.** *Given the approximation in Eq. (5) and  $\langle \Delta\theta_r - \Delta\theta_f, \Delta\theta_n - \Delta\theta_f \rangle \leq 0$ , there exists an  $\alpha < 0$  such that NLS improves GA in unlearning, i.e.,  $\|\theta_r^* - \theta_{f, \text{NLS}}^*\| < \|\theta_r^* - \theta_f^*\|$ , where  $\theta_{f, \text{NLS}}^*$  is the optimal parameters unlearned using GA and NLS, and  $\langle \cdot, \cdot \rangle$  the inner product of two vectors.*

Now we explain the above theorem intuitively. Vector  $\Delta\theta_f - \Delta\theta_r$  is the resultant of Newton’s direction of learning and unlearning. Vector  $\Delta\theta_f - \Delta\theta_n$  is the resultant of Newton’s direction of learning non-target labels and unlearning the target label. When the condition  $\langle \Delta\theta_r - \Delta\theta_f, \Delta\theta_n - \Delta\theta_f \rangle \leq 0$  holds,  $\Delta\theta_n - \Delta\theta_f$  captures the effects of the smoothing term in the unlearning process. If we assume that the exact MU model is able to fully unlearn an example, vector  $\Delta\theta_n$  contributes a direction that pushes the model closer to the exact MU state by leading the model to give the wrong prediction. The illustration of  $\langle \Delta\theta_r - \Delta\theta_f, \Delta\theta_n - \Delta\theta_f \rangle$  is shown in Figure 6 in the Appendix.

The effect of the smoothed term in gradient ascent (GA) with NLS is equivalent to performing a gradient descent optimization with traditionally defined (positive) LS. The gradient of the smoothed term is exactly the same as  $\alpha/K \cdot \sum_{y' \in \mathcal{Y} \setminus y} \nabla \ell(h_{\theta}, (x, y'))$  in both cases.

### 3.4 LABEL SMOOTHING HELPS LOCAL DIFFERENTIAL PRIVACY

When  $\alpha < 0$ , the smoothing term will incur a positive effect in the gradient ascent (GA) step. Label smoothing can also be viewed through the lens of privacy protection. This interpretation stems from the fact that label smoothing reduces the likelihood of a specific label, thereby allowing it to better blend in with other candidate labels. Particularly, we consider a local differential privacy (LDP) guarantee for labels as follows.

**Definition 1 (Label-LDP).** A privacy protection mechanism  $\mathcal{M}$  satisfies  $\epsilon$ -Label-LDP, if for any labels  $y, y', y^{\text{pred}} \in \mathcal{Y}$ ,  $\frac{\mathbb{P}(\mathcal{M}(y)=y^{\text{pred}})}{\mathbb{P}(\mathcal{M}(y')=y^{\text{pred}})} \leq e^\epsilon$ .

The operational meaning of  $\mathcal{M}$  is to guarantee any two labels  $y$  and  $y'$  in the label space, after privatization, have a similar likelihood to become any  $y^{\text{pred}}$  in the label space. That is, the prediction on the forgetting dataset should be similar no matter what the ground-truth label is. The similarity is measured by the privacy budget  $\epsilon \in [0, +\infty)$ . Smaller  $\epsilon$  implies stronger indistinguishability between  $y$  and  $y'$ , and hence, stricter privacy.

Recall  $R_{tr}(\theta) = \sum_{z^{tr} \in D_{tr}} \ell(h_\theta, z^{tr})$ . Denote by  $R_f^{\text{NLS}}(\theta; \alpha) = \sum_{z^{\text{LS}, \alpha} \in D_f} \ell(h_\theta, z^{\text{LS}, \alpha})$ ,  $\alpha < 0$  the empirical risk of forgetting data with NLS. After MU with label smoothing on  $D_f$  by GA, the resulting model can be seen as minimizing the risk  $\gamma_1 \cdot R_{tr}(\theta) - \gamma_2 \cdot R_f^{\text{NLS}}(\theta; \alpha)$ , which is a weighted combination of the risk from two phases: 1) machine learning on  $D_{tr}$  with weight  $\gamma_1 > 0$  and 2) machine unlearning on  $D_f$  with weight  $\gamma_2 > 0$ . By analyzing the risk, we have the following theorem to show NLS in MU induces  $\epsilon$ -Label-LDP for the forgetting data.

**Theorem 3.** Suppose  $\gamma_1 - \gamma_2(1 + \frac{1-K}{K}\alpha) > 0$ . MU using GA+NLS achieves  $\epsilon$ -Label-LDP on  $D_f$  where

$$\epsilon = \left| \log \left( \frac{K}{\alpha} \left( 1 - \frac{\gamma_1}{\gamma_2} \right) + 1 - K \right) \right|, \alpha < 0.$$

Intuitively, when  $\alpha$  is more negative, the privacy of the labels in the forgetting dataset is better. When  $\alpha \rightarrow (1 - \gamma_1/\gamma_2)$ , we have  $\epsilon \rightarrow 0$ , indicating the best label-LDP result, which is the goal of MU. The theorem also warns that  $\alpha$  cannot be arbitrarily negative.

## 4 UGRADSL: A PLUG-AND-PLAY AND GRADIENT-MIXED MU METHOD

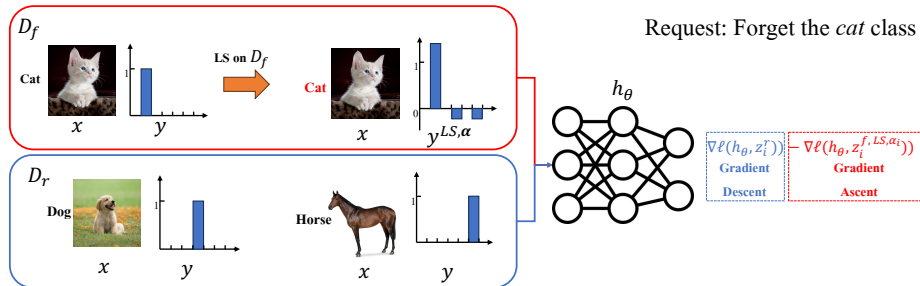


Figure 2: The framework of UGradSL. When there is an unlearning request, we can split the  $D_{tr}$  into  $D_f$  and  $D_r$ . We first apply label smoothing on  $z_i^f = \{x, y\} \in D_f$  to get  $z_i^{\text{LS}, \alpha} = \{x, y^{\text{LS}, \alpha}\}$ , where the smooth rate can be pre-defined or self-adaptive. In the backpropagation process, we apply gradient descent on the data  $z_i^r \in D_r$  and gradient ascent on the smoothed forgetting data  $D_f$ , which is the mix-gradient way.

Given the effect of label smoothing on MU and LDP, we propose our method here. Compared with retraining, Fine-Tune (FT) and GA are much more efficient as illustrated in Section 5 with comparable or better MU performance. FT and GA focus on different perspectives of MU. FT is to transfer the knowledge of the model from  $D_{tr}$  to  $D_r$  using gradient descent (GD) while GA is to remove the knowledge of  $D_f$  from the model.

Due to the flexibility of label smoothing, our method is suitable for the gradient-based methods including FT and GA, making our method a plug-and-play algorithm. **UGradSL is based on GA while UGradSL+ is on FT. Compared with UGradSL, UGradSL+ will lead to a more comprehensive result but with a larger computation cost.**

How to choose the smooth rate  $\alpha$  is worth discussing. Normally, the  $\alpha_i \in \alpha$  for every data point  $z_i^f \in D_f$  can be set identically. To gain better performance, we improve UGradSL and UGradSL+ by taking every data point into consideration and assigning  $\alpha_i$  individually and adaptively based on the distance  $d(h_{\theta_f}(z_i^r), h_{\theta_f}(z_i^f)) \in [0, 1]$  for each  $(z_i^r, z_i^f)$  pair. The intuition is that if an instance  $z_i^f$  resides in a dense neighborhood of  $D_r$ , its inherent deniability is higher and therefore the requirement for “forgetting” is lower and should be reflected through a smaller  $\alpha_i$ . The algorithm is presented in Algorithm 1 and the framework is illustrated in Figure 2. We leave the details of the implementation, complexity analysis and the additional classification results in Appendix D.

Assuming the amount of retained data is significantly larger than the amount of data to be forgotten ( $|D_r| > |D_f|$ ),  $D_f$  will be iterated several times when  $D_r$  is fully iterated once. We calculate the loss using a gradient-mixed method as:

$$L(h_{\theta}, B_f^{\text{NLS}, \alpha}, B_r, p) = p \cdot \sum_{z^r \in B_r} \ell(h_{\theta}, z^r) - (1 - p) \cdot \sum_{z_i^f, \text{NLS}, \alpha_i \in B_f^{\text{NLS}, \alpha}} \ell(h_{\theta}, z_i^f, \text{NLS}, \alpha_i) \quad (8)$$

where  $p \in [0, 1]$  is used to balance GD and GA and the minus sign between two elements on the RHS stands for the GA.  $\alpha$  is the vector for the smoothing rate of every data point  $z_i^f$ .  $h_{\theta}$  is updated according to  $L$  in Eq. (8). UGradSL is similar to UGradSL+ and the dataset used is given in bracket in Algorithm 1. The difference between UGradSL and UGradSL+ is the convergence standard. UGradSL is based on the convergence of  $D_f$  while UGradSL+ is based on  $D_r$ . It should be noted that the Hessian matrix in Theorem 1 is only used in the theoretical proof. In the practical calculation, **there is no need to calculate the Hessian matrix**. Thus, our method does not incur substantially more computation but improves the MU performance on a large scale. We present empirical evidence in Section 5. Compared with applying the label smoothing evenly, the improved version takes the similarity of the data points between  $D_r$  and  $D_f$  into consideration and provides self-adaptive smoothed labels for individual  $z_i^f$  as well as protects label-LDP.

## 5 EXPERIMENTS AND RESULTS

### 5.1 EXPERIMENT SETUP

**Dataset and Model Selection** We validate our method using various datasets in different scales and modalities, including CIFAR-10, CIFAR-100 (Krizhevsky et al., 2009), SVHN (Netzer et al., 2011), CelebA (Liu et al., 2015), Tiny-ImageNet, ImageNet (Deng et al., 2009) and 20 Newsgroups (Lang, 1995) datasets. For the vision and language dataset, we use ResNet-18 (He et al., 2016) and Bert (Devlin et al., 2018) as the backbone model, respectively. Due to the page limit, the details of the training parameter and the additional results of different models including VGG-16 (Simonyan & Zisserman, 2014) and vision transformer (ViT) (Dosovitskiy et al., 2020) are given in the Appendix E.5.

**Baseline Methods** We compare the proposed methods with a series of baselines, including retrain, fine-tuning (FT) (Warnecke et al., 2021; Golatkar et al., 2020), gradient ascent (GA) (Graves et al., 2021; Thudi et al., 2021), unlearning based on the influence function (IU) (Izzo et al., 2021; Koh & Liang, 2017), boundary unlearning (BU) (Chen et al., 2023),  $\ell_1$ -sparse (Jia et al., 2023), random label (RL) (Hayase et al., 2020), SCRUB (Kurmanji et al., 2023), SalUN (Fan et al., 2023), EU- $k$ , CF- $k$  (Goel et al., 2022), GLI (Choi et al., 2024) and PABI (Koloskova et al., 2025). The implementation details of these baselines are given in Appendix E.1.

**Evaluation Metrics** The evaluation metrics we use follow (Jia et al., 2023), where we jointly consider unlearning accuracy (UA), membership inference attack (MIA), remaining accuracy (RA), testing accuracy (TA), and run-time efficiency (RTE). UA is the ratio of incorrect prediction on  $D_f$ , showing the MU performance. TA is the accuracy used to evaluate the performance on the whole testing set  $D_{te}$ , except for the class-wise forgetting because the task is to forget the specific class.

---

**Algorithm 1** UGradSL+: A plug-and-play, efficient, gradient-based MU method using LS. UGradSL can be specified by imposing the [dataset replacement](#) in brackets. If  $\alpha$  is not given, the algorithm turns to the self-adaptive version.

---

**Require:** An almost-converged model  $h_{\theta_{tr}}$  trained with  $D_{tr}$ . The retained dataset  $D_r$ . The forgetting dataset  $D_f$ . Unlearning epochs  $E$ . GA ratio  $p$ . Distance threshold  $\beta$ . The optional smoothing ratio  $\alpha$ .

**Ensure:** The unlearned model  $h_{\theta_f}$ .

```

1: Set the current epoch index as  $t_c \leftarrow 1$ 
2: while  $t_c < E$  do
3:   while  $D_r(D_f)$  is not fully iterated do
4:     Sample a batch  $B_r$  in  $D_r$ 
5:     Sample a batch  $B_f$  from  $D_f$  where  $|B_f| = |B_r|$ 
6:     if  $\alpha$  is not given then ▷ The improved and self-adaptive version
7:       Extract the feature of  $z_i^r$  and  $z_j^f$  using  $h_{\theta_f}$ .
8:       Calculate the distance  $d(h_{\theta_f}(z_i^r), h_{\theta_f}(z_j^f))$  for each  $(z_i^r, z_j^f)$  pair where  $z_i^r \in B_r$ 
       and  $z_j^f \in B_f$ .
9:       For each  $z_j^f$ , count the number  $c_j^f$  of  $z_i^r$  whose  $d(h_{\theta_f}(z_i^r), h_{\theta_f}(z_j^f)) < \beta$ 
10:      Calculate the smooth rate  $\alpha_j = c_j^f / |B_f|$  for each  $z_j^f \in B_f$ 
11:     end if
12:     Update the model using  $B_r, B_f, p$  and  $\alpha_i$  according to Eq. (8)
13:   end while
14:    $t_c \leftarrow t_c + 1$ 
15: end while

```

---

RA is the accuracy on  $D_r$ . To evaluate the effectiveness of "forgetting", we resort to the MIA metrics described in (Jia et al., 2023; Fan et al., 2023), i.e. accuracy of an attack model against target model  $\theta_u$ , such that the score is reported as true negative rate (TNR) on the forget set. Formally, this is a global MIA score (Yeom et al., 2018), which we rewrite as  $\text{MIA}_{\text{Score}} = 1 - \Pr(x_f | \theta_*)$ , where  $x_f \in D_f$  are the forget samples and  $\theta_*$  is the model under test. Overall, we use *Avg. Gap* to quantify the mean performance gap between each unlearning method and the retrained model across all individual metrics above. A lower value indicates better performance.

**Unlearning Paradigm** We mainly consider three unlearning paradigms, including *class-wise forgetting*, *random forgetting*, and *group forgetting*. Class-wise forgetting is to unlearn the whole specific class where we remove one class in  $D_r$  and the corresponding class in  $D_{te}$  completely. Random forgetting across all classes is to unlearn data points belonging to all classes. As a special case of random forgetting, *group forgetting* means that the model is trained to unlearn the group or sub-class of the corresponding super-classes. A more detailed description is given in Appendix E.2.

## 5.2 EXPERIMENT RESULTS

### 5.2.1 CLASS-WISE FORGETTING

We select the class randomly and run class-wise forgetting on five datasets. We report the results of CIFAR-100 / ImageNet and CIFAR-10 in Table 1 and 4, respectively. The results of 20 News-Group and SVHN are given in Appendix E.3. As we can see, UGradSL and UGradSL+ can boost the performance of GA and FT, respectively without an increment in RTE or drop in TA and RA, leading to comprehensive satisfaction in the main metrics, despite randomness in  $D_f$ , showing the robustness and flexibility of our methods in MU regardless of the size of the dataset and the data modality. Moreover, in terms of *Avg. Gap*, the proposed method shows its similarity to the retrained model.

### 5.2.2 RANDOM FORGETTING

We select data randomly from every class as  $D_f$ , making sure all the classes are selected and the size of  $D_f$  is 10% of the  $D_{tr}$ . We report the results of CIFAR-100 and TinyImageNet in Table 2.

Table 1: Results of class-wise forgetting in CIFAR-100 and ImageNet. The best comprehensive metrics are **bold**.

Method	CIFAR-100						ImageNet					
	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap (↓)	RTE (↓, min)	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap (↓)	RTE (↓, hr)
Retrain	100.00 <sub>±0.00</sub>	100.00 <sub>±0.00</sub>	99.96 <sub>±0.01</sub>	71.10 <sub>±0.12</sub>	-	26.95	100.00 <sub>±0.00</sub>	100.00 <sub>±0.00</sub>	71.62 <sub>±0.12</sub>	69.57 <sub>±0.07</sub>	-	26.18
FT	0.67 <sub>±0.38</sub>	27.20 <sub>±1.84</sub>	99.96 <sub>±0.01</sub>	71.46 <sub>±0.09</sub>	43.12	1.74	52.42 <sub>±15.81</sub>	55.86 <sub>±18.02</sub>	70.66 <sub>±2.54</sub>	69.25 <sub>±0.78</sub>	23.25	2.87
GA	99.00 <sub>±0.57</sub>	99.07 <sub>±0.50</sub>	77.83 <sub>±2.07</sub>	53.73 <sub>±0.96</sub>	10.36	<b>0.06</b>	81.23 <sub>±0.69</sub>	83.52 <sub>±2.08</sub>	66.00 <sub>±0.03</sub>	64.72 <sub>±0.02</sub>	11.43	<b>0.01</b>
IU	2.07 <sub>±1.65</sub>	33.20 <sub>±8.83</sub>	99.96 <sub>±0.01</sub>	71.39 <sub>±0.19</sub>	41.26	1.24	33.54 <sub>±19.46</sub>	49.83 <sub>±21.57</sub>	66.25 <sub>±1.99</sub>	66.28 <sub>±1.19</sub>	31.32	1.51
BE	99.07 <sub>±0.34</sub>	99.00 <sub>±0.49</sub>	70.81 <sub>±2.69</sub>	49.85 <sub>±1.32</sub>	13.08	0.55	98.62 <sub>±0.58</sub>	0.15 <sub>±0.11</sub>	53.13 <sub>±0.27</sub>	56.72 <sub>±0.31</sub>	33.14	0.24
BS	98.87 <sub>±0.57</sub>	98.73 <sub>±0.68</sub>	71.16 <sub>±2.60</sub>	50.03 <sub>±1.36</sub>	13.06	0.77	98.85 <sub>±0.50</sub>	0.13 <sub>±0.12</sub>	53.35 <sub>±0.16</sub>	56.93 <sub>±0.03</sub>	32.98	0.37
$\ell_1$ -sparse	98.97 <sub>±1.03</sub>	100.00 <sub>±0.00</sub>	86.99 <sub>±0.76</sub>	79.08 <sub>±0.75</sub>	5.50	0.15	100.00 <sub>±0.00</sub>	100.00 <sub>±0.00</sub>	39.01 <sub>±1.03</sub>	44.62 <sub>±0.91</sub>	14.39	0.16
RL	99.80 <sub>±0.35</sub>	100.00 <sub>±0.00</sub>	99.97 <sub>±0.62</sub>	77.31 <sub>±0.35</sub>	1.61	1.10	100.00 <sub>±0.00</sub>	100.00 <sub>±0.00</sub>	62.06 <sub>±4.19</sub>	62.93 <sub>±0.45</sub>	4.05	1.17
EU-k	100.00 <sub>±0.00</sub>	0.00 <sub>±0.00</sub>	63.79 <sub>±1.10</sub>	43.90 <sub>±0.73</sub>	40.84	4.50	100.00 <sub>±0.00</sub>	0.00 <sub>±0.00</sub>	32.99 <sub>±0.07</sub>	37.19 <sub>±0.15</sub>	42.75	0.62
CF-k	100.00 <sub>±0.00</sub>	0.00 <sub>±0.00</sub>	94.88 <sub>±0.46</sub>	61.32 <sub>±1.17</sub>	28.72	3.01	99.79 <sub>±0.36</sub>	0.00 <sub>±0.00</sub>	66.84 <sub>±0.03</sub>	68.35 <sub>±0.28</sub>	26.55	1.25
SCRUB	30.07 <sub>±49.48</sub>	66.90 <sub>±29.19</sub>	99.98 <sub>±0.01</sub>	77.97 <sub>±0.56</sub>	27.56	1.07	56.59 <sub>±2.17</sub>	75.59 <sub>±1.19</sub>	66.98 <sub>±0.11</sub>	68.24 <sub>±0.07</sub>	18.45	0.21
SaUN	99.90 <sub>±0.01</sub>	99.96 <sub>±0.00</sub>	99.98 <sub>±0.01</sub>	75.02 <sub>±0.10</sub>	1.02	2.15	100.00 <sub>±0.00</sub>	100.00 <sub>±0.00</sub>	63.00 <sub>±5.03</sub>	62.72 <sub>±0.31</sub>	3.87	1.95
GLI	39.78 <sub>±6.74</sub>	69.63 <sub>±7.44</sub>	95.57 <sub>±2.11</sub>	69.63 <sub>±0.60</sub>	24.11	1.03	53.38 <sub>±2.96</sub>	73.31 <sub>±3.22</sub>	73.01 <sub>±0.11</sub>	63.23 <sub>±0.06</sub>	20.26	3.79
PABI	100.00 <sub>±0.00</sub>	100.00 <sub>±0.00</sub>	98.94 <sub>±0.16</sub>	73.41 <sub>±0.09</sub>	0.83	20.09	-	-	-	-	-	-
UGradSL	66.59 <sub>±0.90</sub>	90.96 <sub>±5.05</sub>	95.45 <sub>±1.42</sub>	70.34 <sub>±1.78</sub>	11.93	0.07	100.00 <sub>±0.00</sub>	100.00 <sub>±0.00</sub>	76.91 <sub>±1.82</sub>	65.94 <sub>±1.35</sub>	<b>2.23</b>	<b>0.01</b>
UGradSL+	100.00 <sub>±0.00</sub>	100.00 <sub>±0.00</sub>	98.44 <sub>±0.62</sub>	72.12 <sub>±0.70</sub>	<b>0.64</b>	3.37	100.00 <sub>±0.00</sub>	100.00 <sub>±0.00</sub>	78.16 <sub>±0.07</sub>	66.84 <sub>±0.06</sub>	2.32	4.19

Table 2: Results of random forgetting in CIFAR-100 and Tiny-ImageNet. The best comprehensive metrics are **bold**.

Method	CIFAR-100						Tiny-ImageNet					
	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap (↓)	RTE (↓, min)	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap (↓)	RTE (↓, min)
Retrain	29.47 <sub>±1.59</sub>	53.50 <sub>±1.19</sub>	99.98 <sub>±0.01</sub>	70.51 <sub>±1.17</sub>	-	25.01	49.35 <sub>±0.38</sub>	58.44 <sub>±0.89</sub>	83.80 <sub>±0.29</sub>	59.66 <sub>±0.44</sub>	-	235.68
FT	2.55 <sub>±0.03</sub>	10.59 <sub>±0.27</sub>	99.95 <sub>±0.01</sub>	75.95 <sub>±0.05</sub>	18.83	1.95	29.23 <sub>±0.29</sub>	37.02 <sub>±0.33</sub>	82.51 <sub>±0.20</sub>	60.96 <sub>±0.23</sub>	11.03	18.61
GA	2.58 <sub>±0.06</sub>	5.95 <sub>±0.17</sub>	97.45 <sub>±0.02</sub>	76.09 <sub>±0.01</sub>	20.64	0.29	19.34 <sub>±1.67</sub>	25.19 <sub>±0.88</sub>	81.51 <sub>±1.56</sub>	59.66 <sub>±0.61</sub>	16.39	8.65
IU	15.71 <sub>±5.19</sub>	18.69 <sub>±4.12</sub>	84.65 <sub>±5.29</sub>	62.20 <sub>±4.17</sub>	18.05	1.20	60.61 <sub>±0.01</sub>	83.67 <sub>±0.15</sub>	16.36 <sub>±0.37</sub>	23.44 <sub>±0.29</sub>	35.04	7.30
BE	0.01 <sub>±0.00</sub>	1.45 <sub>±0.02</sub>	99.97 <sub>±0.18</sub>	78.26 <sub>±0.00</sub>	22.32	<b>0.24</b>	17.65 <sub>±0.31</sub>	24.48 <sub>±0.42</sub>	82.85 <sub>±0.20</sub>	58.16 <sub>±0.08</sub>	17.03	<b>3.53</b>
BS	2.20 <sub>±1.21</sub>	10.73 <sub>±0.37</sub>	98.22 <sub>±1.26</sub>	70.23 <sub>±1.67</sub>	18.02	0.34	19.47 <sub>±0.69</sub>	25.45 <sub>±0.15</sub>	81.23 <sub>±0.74</sub>	56.75 <sub>±0.80</sub>	17.09	5.63
$\ell_1$ -sparse	8.19 <sub>±0.38</sub>	19.11 <sub>±0.52</sub>	88.39 <sub>±0.31</sub>	80.26 <sub>±0.16</sub>	19.25	1.00	35.73 <sub>±0.35</sub>	41.98 <sub>±0.73</sub>	78.19 <sub>±0.05</sub>	61.44 <sub>±0.12</sub>	9.37	23.40
RL	4.06 <sub>±0.37</sub>	50.12 <sub>±3.48</sub>	99.92 <sub>±0.01</sub>	71.30 <sub>±0.12</sub>	7.41	1.20	40.52 <sub>±0.15</sub>	59.01 <sub>±0.76</sub>	77.58 <sub>±0.06</sub>	60.18 <sub>±0.19</sub>	4.04	27.08
EU-k	1.73 <sub>±0.06</sub>	3.33 <sub>±0.07</sub>	98.44 <sub>±0.05</sub>	59.92 <sub>±0.43</sub>	22.51	1.96	33.55 <sub>±0.35</sub>	22.19 <sub>±1.75</sub>	81.41 <sub>±0.27</sub>	58.08 <sub>±0.21</sub>	14.01	20.02
CF-k	0.07 <sub>±0.02</sub>	0.47 <sub>±0.16</sub>	99.98 <sub>±0.01</sub>	67.86 <sub>±0.12</sub>	21.27	0.88	19.31 <sub>±0.38</sub>	23.22 <sub>±2.28</sub>	81.59 <sub>±0.37</sub>	58.15 <sub>±0.19</sub>	17.25	13.18
SCRUB	0.09 <sub>±0.59</sub>	4.01 <sub>±1.25</sub>	99.97 <sub>±0.34</sub>	77.45 <sub>±0.26</sub>	21.46	1.06	20.11 <sub>±1.15</sub>	25.35 <sub>±7.53</sub>	80.91 <sub>±0.77</sub>	60.11 <sub>±0.99</sub>	16.42	25.79
SaUN	35.23 <sub>±0.32</sub>	89.39 <sub>±0.46</sub>	99.53 <sub>±0.04</sub>	64.26 <sub>±0.58</sub>	12.10	3.33	40.39 <sub>±0.15</sub>	52.32 <sub>±0.67</sub>	77.60 <sub>±0.11</sub>	60.30 <sub>±0.31</sub>	5.48	34.42
GLI	2.88 <sub>±1.51</sub>	9.33 <sub>±2.36</sub>	97.16 <sub>±1.48</sub>	72.04 <sub>±0.30</sub>	18.78	0.63	38.40 <sub>±1.74</sub>	67.87 <sub>±2.20</sub>	98.37 <sub>±0.27</sub>	61.53 <sub>±0.30</sub>	9.21	22.92
PABI	28.33 <sub>±0.74</sub>	39.31 <sub>±0.88</sub>	99.14 <sub>±0.01</sub>	72.00 <sub>±0.20</sub>	4.42	19.10	99.90 <sub>±0.03</sub>	66.46 <sub>±56.83</sub>	0.50 <sub>±0.01</sub>	0.00 <sub>±0.00</sub>	50.38	54.58
UGradSL	18.36 <sub>±0.17</sub>	40.71 <sub>±0.13</sub>	98.38 <sub>±0.03</sub>	68.23 <sub>±0.16</sub>	6.95	0.55	40.73 <sub>±0.71</sub>	37.58 <sub>±0.21</sub>	67.30 <sub>±0.04</sub>	50.38 <sub>±0.77</sub>	13.82	9.47
UGradSL+	21.69 <sub>±0.59</sub>	49.47 <sub>±1.25</sub>	99.87 <sub>±0.34</sub>	73.60 <sub>±0.26</sub>	<b>3.75</b>	3.52	53.06 <sub>±1.27</sub>	59.46 <sub>±1.01</sub>	81.38 <sub>±0.75</sub>	52.52 <sub>±0.84</sub>	<b>3.57</b>	25.93

Compared with class-wise forgetting, it is harder to improve the MU performance and still keep the RA and TA close to the retrained model. Benefit from the mix-gradient design, the proposed method can make a good balance between forgetting  $D_f$  and retaining the knowledge in  $D_r$ . The rest of the experiments are given in Appendix E.4.

### 5.2.3 GROUP FORGETTING

Although group forgetting can be seen as part of random forgetting, we want to highlight its use case here due to its practical impacts on *e.g.*, facial attributes classification. The identities can be regarded as the subgroup in the attributes.

Table 3: Results of Group Forgetting on CIFAR-20 and CelebA. For CIFAR-20, the model is trained to classify 20 super-classes, with  $D_f$  representing one of five subclasses within a single super-class. In the CelebA dataset, the model performs binary classification to determine whether a person is smiling, with  $D_f$  selected based on specific identities. The best comprehensive metrics are **bold**.

Method	CIFAR-20						CelebA					
	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap (↓)	RTE (↓, min)	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap (↓)	RTE (↓, min)
Retrain	13.33 <sub>±1.64</sub>	28.47 <sub>±0.75</sub>	99.94 <sub>±0.01</sub>	81.23 <sub>±0.13</sub>	-	27.35	6.74 <sub>±0.26</sub>	9.77 <sub>±1.49</sub>	94.38 <sub>±0.49</sub>	91.78 <sub>±0.33</sub>	-	258.69
FT	1.00 <sub>±0.43</sub>	2.73 <sub>±0.52</sub>	99.37 <sub>±0.08</sub>	79.02 <sub>±0.03</sub>	10.21	7.47	5.36 <sub>±0.17</sub>	5.87 <sub>±0.11</sub>	93.91 <sub>±0.04</sub>	93.18 <sub>±0.03</sub>	1.79	25.94
GA	87.93 <sub>±2.92</sub>	88.93 <sub>±2.33</sub>	81.46 <sub>±0.77</sub>	64.07 <sub>±0.95</sub>	42.68	<b>0.11</b>	6.00 <sub>±0.16</sub>	5.76 <sub>±0.14</sub>	92.86 <sub>±0.13</sub>	92.52 <sub>±0.08</sub>	1.75	1.20
IU	0.00 <sub>±0.00</sub>	2.07 <sub>±1.29</sub>	99.95 <sub>±0.01</sub>	80.92 <sub>±0.34</sub>	10.01	1.10	5.90 <sub>±0.11</sub>	4.91 <sub>±0.30</sub>	93.05 <sub>±0.01</sub>	92.62 <sub>±0.01</sub>	1.97	219.77
BE	89.07 <sub>±1.39</sub>	91.77 <sub>±1.75</sub>	76.36 <sub>±0.92</sub>	60.17 <sub>±0.92</sub>	45.91	0.33	11.50 <sub>±0.80</sub>	48.41 <sub>±8.86</sub>	88.37 <sub>±0.81</sub>	88.07 <sub>±0.81</sub>	13.28	48.91
BS	88.60 <sub>±1.13</sub>	90.67 <sub>±1.18</sub>	76.70 <sub>±1.08</sub>	60.41 <sub>±1.17</sub>	45.38	0.29	8.95 <sub>±5.11</sub>	27.35 <sub>±30.20</sub>	91.00 <sub>±0.22</sub>	90.63 <sub>±5.65</sub>	6.08	50.99
$\ell_1$ -sparse	0.13 <sub>±0.09</sub>	2.27 <sub>±0.57</sub>	99.57 <sub>±1.08</sub>	80.44 <sub>±0.08</sub>	10.14	0.38	9.46 <sub>±1.82</sub>	36.91 <sub>±30.96</sub>	90.52 <sub>±5.22</sub>	90.35 <sub>±1.77</sub>	8.79	37.49
RL	56.93 <sub>±3.24</sub>	98.60 <sub>±0.29</sub>	99.92 <sub>±0.01</sub>	80.28 <sub>±0.05</sub>	28.67	0.37	8.31 <sub>±0.43</sub>	28.55 <sub>±16.74</sub>	91.85 <sub>±0.51</sub>	91.62 <sub>±0.42</sub>	5.76	40.09
EU-k	8.00 <sub>±1.57</sub>	16.33 <sub>±7.18</sub>	97.07 <sub>±0.18</sub>	69.67 <sub>±0.35</sub>	7.98	0.87	7.20 <sub>±0.19</sub>	18.77 <sub>±3.69</sub>	92.55 <sub>±0.30</sub>	91.04 <sub>±0.67</sub>	3.01	1.98
CF-k	0.00 <sub>±0.00</sub>	0.80 <sub>±0.40</sub>	99.98 <sub>±0.01</sub>	77.46 <sub>±0.03</sub>	11.20	1.21	5.46 <sub>±0.32</sub>	17.26 <sub>±0.08</sub>	94.45 <sub>±0.04</sub>	92.72 <sub>±0.04</sub>	2.45	1.60
SCRUB	0.00 <sub>±0.00</sub>	1.13 <sub>±0.34</sub>	99.93 <sub>±0.01</sub>	81.05 <sub>±0.20</sub>	10.21	0.30	8.78 <sub>±0.77</sub>	13.37 <sub>±5.22</sub>	91.21 <sub>±0.86</sub>	90.65 <sub>±0.86</sub>	2.49	70.13
SaUN	52.93 <sub>±2.21</sub>	99.80 <sub>±0.35</sub>	99.55 <sub>±0.00</sub>	76.48 <sub>±0.26</sub>	29.02	2.88	6.53 <sub>±0.28</sub>	25.57 <sub>±8.22</sub>	92.97 <sub>±0.03</sub>	92.27 <sub>±0.07</sub>	4.48	83.43
GLI	22.22 <sub>±2.47</sub>	34.26 <sub>±4.88</sub>	91.12 <sub>±0.66</sub>	75.13 <sub>±0.73</sub>	7.40	1.01	6.43 <sub>±0.43</sub>	5.70 <sub>±0.74</sub>	92.59 <sub>±0.08</sub>	92.36 <sub>±0.11</sub>	1.69	<b>0.52</b>
PABI	82.96 <sub>±9.83</sub>	90.83 <sub>±9.88</sub>	99.29 <sub>±0.47</sub>	81.47 <sub>±0.02</sub>	33.22	19.96	3.71 <sub>±0.19</sub>	45.16 <sub>±1.97</sub>	99.83 <sub>±0.01</sub>	99.57 <sub>±0.00</sub>	12.92	38.87
UGradSL	22.87 <sub>±0.90</sub>	38.93 <sub>±1.57</sub>	97.20 <sub>±0.19</sub>	75.84 <sub>±0.16</sub>	<b>7.03</b>	0.13	6.29 <sub>±1.41</sub>	5.73 <sub>±3.50</sub>	93.44 <sub>±0.14</sub>	92.80 <sub>±0.27</sub>	1.61	2.17
UGradSL+	78.44 <sub>±1.19</sub>	88.67 <sub>±0.35</sub>	97.93 <sub>±0.71</sub>	79.77 <sub>±0.58</sub>	32.20	8.12	6.88 <sub>±0.63</sub>	8.88 <sub>±0.52</sub>	93.39 <sub>±0.41</sub>	92.99 <sub>±0.12</sub>	<b>0.81</b>	16.78

**CIFAR-10 and CIFAR-100** share the same image dataset while CIFAR-100 is labeled with 100 fine-grained classes and 20 coarse (super) classes (Krizhevsky et al., 2009; Chundawat et al., 2023). We train a model to classify 20 super classes using CIFAR-100 training set. The setting of the *group forgetting within one coarse class* is to remove one fine-grained class from one super class in CIFAR-100 datasets. For example, there are five fine-grained fishes in the *Fish* coarse class and we want to remove one fine-grained fish from the model. Different from class-wise forgetting, we do not modify the testing set. We report the group forgetting in Table 3.

**CelebA** We select CelebA dataset as another real-world case and show the results in Table 3. We train a binary classification model to classify whether the person is smile or not. There are 8192 identities in the training set and we select 1% of the identities (82 identities) as  $D_f$ . Both smiling and non-smiling images are in  $D_f$ . This experiment has significant practical meaning, since the biometric, such as identity and fingerprint, needs more privacy protection (Minaee et al., 2023). Compared with baseline methods, our method can forget the identity information better without forgetting too much remaining information in the dataset. This paradigm provides a practical usage of MU and our methods provide a faster and more reliable way to improve the MU performance.

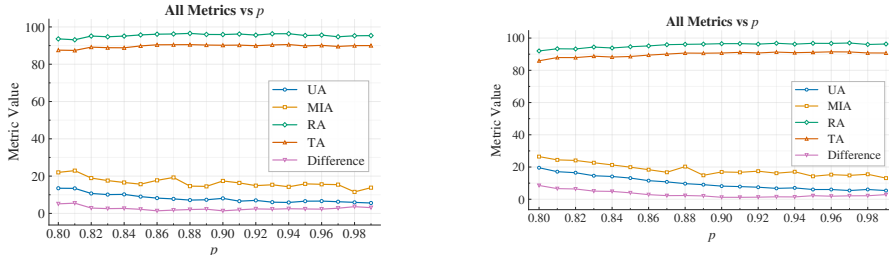


Figure 3: Ablation study on the GA ratio  $p$  for random forgetting (10%) on CIFAR-10. Our methods (left: UGradSL, right: UGradSL+) remain relatively stable across a wide range of  $p$  under multiple evaluation metrics.

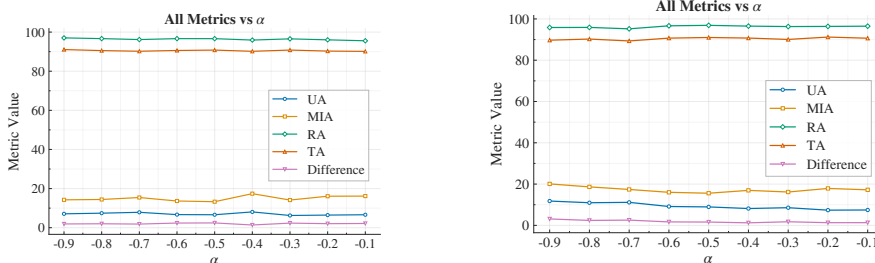


Figure 4: Ablation study on the smoothing rate  $\alpha$  for random forgetting (10%) on CIFAR-10. Our methods (left: UGradSL, right: UGradSL+) remain relatively stable across a wide range of  $\alpha$  under multiple evaluation metrics.

### 5.3 ABLATION STUDIES

To evaluate the robustness of our method, we conduct ablation studies on the forgetting set size, the GA ratio  $p$ , and the smoothing rate  $\alpha$ . For the forgetting set size, we do the experiments on CIFAR-10 and CIFAR-100. The results in Tables 17 and 18 show that our method is consistently robust to the size of  $D_f$  and always outperforms the other baselines. For the trade-off parameter  $p$ , Table 19 indicates that pure GA is unstable, and its performance becomes more stable as more GA steps are combined with our method. In practice, we therefore choose  $p > 0.9$  to stabilize performance while still leveraging the effect of GA for unlearning. We further study random forgetting (10%) on CIFAR-10, fixing  $\alpha = -0.4$  for both UGradSL and UGradSL+, and sweep  $p$  in  $[0.8, 1.0]$  with a step size of 0.01. The resulting UA, MIA, RA, TA, and Avg. Gap curves in Figure 3 show that our methods are relatively stable with respect to  $p$ . For the smoothing rate  $\alpha$ , we fix  $p = 0.9$  and vary  $\alpha$  from  $-0.9$  to  $-0.1$  with a step size of 0.1. The results in Figure 4 demonstrate that both

UGradSL and UGradSL+ remain stable. Overall, these ablations show that our methods are robust to the choice of hyperparameters.

Table 4: Results of class-wise forgetting and random forgetting on CIFAR-10 with additional (add.) MIA. The best comprehensive metrics are **bold**. Avg. Gap is calculated with additional MIA.

	Class-wise								Random							
	UA	MIA <sub>Score</sub>	RA	TA	Add. MIA	Avg. Gap (↓)	RTE (↓, min)	UA	MIA <sub>Score</sub>	RA	TA	Add. MIA	Avg. Gap (↓)	RTE (↓, min)		
Retrain	100.00 <sub>±0.00</sub>	100.00 <sub>±0.00</sub>	98.19 <sub>±1.14</sub>	94.50 <sub>±0.34</sub>	99.23 <sub>±0.08</sub>	-	24.62	8.07 <sub>±0.47</sub>	17.41 <sub>±0.69</sub>	100.00 <sub>±0.01</sub>	91.61 <sub>±0.24</sub>	50.69 <sub>±0.73</sub>	-	24.66		
FT	22.71 <sub>±5.31</sub>	79.21 <sub>±8.60</sub>	99.82 <sub>±0.09</sub>	94.13 <sub>±0.14</sub>	99.09 <sub>±0.07</sub>	<b>20.04</b>	2.02	1.10 <sub>±0.19</sub>	4.06 <sub>±0.41</sub>	99.83 <sub>±0.03</sub>	93.70 <sub>±0.10</sub>	54.05 <sub>±0.21</sub>	5.19	1.58		
GA	25.19 <sub>±11.38</sub>	73.48 <sub>±9.68</sub>	96.84 <sub>±0.58</sub>	73.10 <sub>±1.62</sub>	99.43 <sub>±0.09</sub>	24.86	<b>0.08</b>	0.56 <sub>±0.01</sub>	1.19 <sub>±0.05</sub>	99.48 <sub>±0.02</sub>	94.55 <sub>±0.05</sub>	55.04 <sub>±0.66</sub>	6.31	<b>0.31</b>		
IU	83.92 <sub>±1.16</sub>	92.59 <sub>±1.41</sub>	98.77 <sub>±0.12</sub>	92.64 <sub>±0.23</sub>	99.71 <sub>±0.07</sub>	5.28	1.18	17.51 <sub>±2.19</sub>	21.39 <sub>±1.70</sub>	83.28 <sub>±2.44</sub>	78.13 <sub>±2.85</sub>	53.98 <sub>±0.55</sub>	9.37	1.18		
BE	64.93 <sub>±0.01</sub>	98.19 <sub>±0.00</sub>	99.47 <sub>±0.00</sub>	94.00 <sub>±0.11</sub>	99.60 <sub>±0.02</sub>	7.81	0.20	0.00 <sub>±0.00</sub>	0.26 <sub>±0.02</sub>	100.00 <sub>±0.00</sub>	95.35 <sub>±0.18</sub>	55.41 <sub>±0.49</sub>	6.74	3.17		
BS	93.69 <sub>±4.32</sub>	99.82 <sub>±0.04</sub>	97.69 <sub>±1.29</sub>	92.89 <sub>±1.26</sub>	99.56 <sub>±0.10</sub>	1.79	0.29	0.48 <sub>±0.07</sub>	1.16 <sub>±0.04</sub>	99.47 <sub>±0.01</sub>	94.58 <sub>±0.03</sub>	55.88 <sub>±0.72</sub>	6.51	1.41		
$\ell_1$ -sparse	100.00 <sub>±0.00</sub>	100.00 <sub>±0.00</sub>	97.86 <sub>±1.29</sub>	96.11 <sub>±1.26</sub>	99.02 <sub>±0.15</sub>	0.43	1.00	2.80 <sub>±0.37</sub>	18.59 <sub>±3.48</sub>	99.97 <sub>±0.01</sub>	94.08 <sub>±0.12</sub>	55.01 <sub>±0.49</sub>	2.65	1.98		
RL	99.99 <sub>±0.01</sub>	100.00 <sub>±0.00</sub>	100.00 <sub>±0.00</sub>	95.50 <sub>±1.11</sub>	99.08 <sub>±0.07</sub>	0.59	1.04	2.27 <sub>±0.29</sub>	14.37 <sub>±1.01</sub>	99.98 <sub>±0.01</sub>	94.14 <sub>±0.18</sub>	52.17 <sub>±0.87</sub>	2.57	1.98		
EU- $k$	100.00 <sub>±0.00</sub>	100.00 <sub>±0.00</sub>	100.00 <sub>±0.00</sub>	75.04 <sub>±1.10</sub>	99.89 <sub>±0.18</sub>	4.39	1.45	0.00 <sub>±0.00</sub>	0.50 <sub>±0.30</sub>	99.99 <sub>±0.01</sub>	77.21 <sub>±1.21</sub>	61.88 <sub>±1.33</sub>	10.12	1.58		
CF- $k$	100.00 <sub>±0.00</sub>	100.00 <sub>±0.00</sub>	100.00 <sub>±0.00</sub>	78.95 <sub>±0.53</sub>	100.00 <sub>±0.00</sub>	3.63	1.32	0.00 <sub>±0.00</sub>	0.00 <sub>±0.00</sub>	100.00 <sub>±0.00</sub>	80.98 <sub>±0.27</sub>	69.91 <sub>±1.33</sub>	11.07	1.47		
SCRUB	100.00 <sub>±0.00</sub>	100.00 <sub>±0.00</sub>	99.93 <sub>±0.01</sub>	95.22 <sub>±0.07</sub>	100.00 <sub>±0.00</sub>	0.65	1.09	0.70 <sub>±0.59</sub>	3.88 <sub>±1.25</sub>	99.59 <sub>±0.34</sub>	94.22 <sub>±0.26</sub>	55.33 <sub>±0.59</sub>	5.71	4.05		
SalIN	90.74 <sub>±13.01</sub>	100.00 <sub>±0.00</sub>	98.29 <sub>±0.34</sub>	80.49 <sub>±1.21</sub>	98.63 <sub>±0.59</sub>	4.78	2.22	46.95 <sub>±0.15</sub>	86.33 <sub>±1.29</sub>	97.75 <sub>±0.42</sub>	77.22 <sub>±0.77</sub>	69.95 <sub>±0.12</sub>	28.74	2.42		
UGradSL	94.99 <sub>±4.35</sub>	97.95 <sub>±1.78</sub>	95.47 <sub>±4.08</sub>	86.78 <sub>±5.68</sub>	99.94 <sub>±0.01</sub>	3.64	0.22	5.87 <sub>±0.51</sub>	13.33 <sub>±0.70</sub>	98.82 <sub>±0.28</sub>	92.17 <sub>±0.23</sub>	53.54 <sub>±0.97</sub>	<b>2.17</b>	0.45		
UGradSL+	100.00 <sub>±0.00</sub>	100.00 <sub>±0.00</sub>	99.29 <sub>±0.01</sub>	94.29 <sub>±0.07</sub>	100.00 <sub>±0.00</sub>	<b>0.41</b>	3.07	6.03 <sub>±0.17</sub>	10.65 <sub>±0.13</sub>	99.79 <sub>±0.03</sub>	93.64 <sub>±0.16</sub>	52.29 <sub>±0.85</sub>	2.53	3.07		

## 5.4 DISCUSSION

**Influence Function in Deep Learning** Influence functions were originally proposed for the convex function. As given in Section 3.2, we apply the influence function to the converged model, which can be regarded as a local convex model. A plot of loss landscape of the retrained model  $\theta_r$  on CIFAR-10 dataset is given in Figure 7 in Appendix.

**MIA as a Proxy for "Forgetfulness".** Given a model  $\theta_*$ , we can evaluate the degree of its generalization by running a membership inference attack on the model. In the context of the current work, generalization is equivalent to the degree of "forgetfulness" that the forgetting algorithm achieves. Given the distribution of model response observations  $A_f = \mathcal{A}(\theta_*, \mathcal{D}_f)$  and  $A_{te} = \mathcal{A}(\theta_*, \mathcal{D}_{te})$ , where  $\mathcal{A}$  is an adversary and  $A = A_f \cup A_{te}$  is the observation visible to  $\mathcal{A}$ , one can get the degree of generalization by analyzing the observations. In the context of MU, the most straightforward way is to get the accuracy of  $\mathcal{A}$  on the seen and unseen samples ( $\mathcal{D}_{te}$  and  $\mathcal{D}_f$  respectively). This could be done by computing the  $(TP + TN) / (|D_f| + |D_{te}|)$ , where the true positive (TP) predictions correspond to "seen" samples, and true negative (TN) predictions are "unseen" samples. We conducted the experiments on CIFAR-10 both for class-wise and random forgetting. The results are given in Table 4, where Avg. Gap is calculated with additional MIA. We assume that the distribution of  $D_{tr}$  and  $D_{te}$  should be the same. For class-wise forgetting, the additional MIA is almost 1 because  $D_f$  is a separate single class and the distribution of  $D_f$  and  $D_{te}$  without the corresponding class are totally different. For random forgetting, the additional MIA is almost 0.5 because  $D_f$  is randomly selected from  $D_{tr}$  and the distribution of  $D_f$  and  $D_{te}$  should be the same. The plots of loss distribution for random and class-wise forgetting are given in Figure 8 in the Appendix. In Table 4, the proposed methods still outperform the other baseline methods, showing the robustness to the other MIA auditing methods and the generalization capability in privacy preservation.

**Difference between UGradSL and UGradSL+** Although two methods are similar in the mathematical formulation, there exists fundamental difference in their design and behavior. Compared with UGradSL, UGradSL+ can be more stable and less sensitive due to its origin from FT. As shown in the experiment results in the tables, UGradSL+ can always perform as top-tier methods. However, the RTE of UGradSL+ would be higher. We present more analysis in Appendix E.11.

The study of the Streisand effect (Jansen & Martin, 2015) and gradient analysis are given in Appendix E.6 and E.10, respectively.

## 6 CONCLUSIONS

We have proposed UGradSL, a plug-and-play, efficient, gradient-based MU method using smoothed labels. Theoretical proofs and extensive numerical experiments have demonstrated the effectiveness of the proposed method.

## 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 D and E. We also provide pseudo-code in Algorithm 1 and will release the code upon acceptance.

## ACKNOWLEDGMENTS

Z. Di and Y. Liu are partially supported by the National Science Foundation (NSF) under grants IIS-2143895 and IIS-2040800. The work of J. Jia and S. Liu is supported in part by the Open Philanthropy Research Award.

## REFERENCES

- Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*, pp. 308–318, 2016.
- Samyadeep Basu, Xuchen You, and Soheil Feizi. On second-order group influence functions for black-box predictions. In *International Conference on Machine Learning*, pp. 715–724. PMLR, 2020.
- Alexander Becker and Thomas Liebig. Evaluating machine unlearning via epistemic uncertainty. *arXiv preprint arXiv:2208.10836*, 2022.
- Lucas Bourtole, Varun Chandrasekaran, Christopher A Choquette-Choo, Hengrui Jia, Adelin Travers, Baiwu Zhang, David Lie, and Nicolas Papernot. Machine unlearning. In *2021 IEEE Symposium on Security and Privacy (SP)*, pp. 141–159. IEEE, 2021.
- Marc-Etienne Brunet, Colleen Alkalay-Houlihan, Ashton Anderson, and Richard Zemel. Understanding the origins of bias in word embeddings. In *International conference on machine learning*, pp. 803–811. PMLR, 2019.
- 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.
- Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. Extracting training data from large language models. In *30th USENIX security symposium (USENIX Security 21)*, pp. 2633–2650, 2021.
- Min Chen, Zhikun Zhang, Tianhao Wang, Michael Backes, Mathias Humbert, and Yang Zhang. When machine unlearning jeopardizes privacy. In *Proceedings of the 2021 ACM SIGSAC conference on computer and communications security*, pp. 896–911, 2021a.
- Min Chen, Zhikun Zhang, Tianhao Wang, Michael Backes, Mathias Humbert, and Yang Zhang. When machine unlearning jeopardizes privacy. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*, pp. 896–911, 2021b.
- Min Chen, Zhikun Zhang, Tianhao Wang, Michael Backes, Mathias Humbert, and Yang Zhang. Graph unlearning. In *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security*, pp. 499–513, 2022.
- Min Chen, Weizhuo Gao, Gaoyang Liu, Kai Peng, and Chen Wang. Boundary unlearning. *arXiv preprint arXiv:2303.11570*, 2023.

- Eli Chien, Chao Pan, and Olgica Milenkovic. Certified graph unlearning. *arXiv preprint arXiv:2206.09140*, 2022.
- Dasol Choi, Soora Choi, Eunsun Lee, Jinwoo Seo, and Dongbin Na. Towards efficient machine unlearning with data augmentation: Guided loss-increasing (gli) to prevent the catastrophic model utility drop. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 93–102, 2024.
- Christopher A Choquette-Choo, Florian Tramer, Nicholas Carlini, and Nicolas Papernot. Label-only membership inference attacks. In *International conference on machine learning*, pp. 1964–1974. PMLR, 2021.
- Jan Chorowski and Navdeep Jaitly. Towards better decoding and language model integration in sequence to sequence models. *arXiv preprint arXiv:1612.02695*, 2016.
- Vikram S Chundawat, Ayush K Tarun, Murari Mandal, and Mohan Kankanhalli. Can bad teaching induce forgetting? unlearning in deep networks using an incompetent teacher. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 7210–7217, 2023.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Jimmy Z Di, Jack Douglas, Jayadev Acharya, Gautam Kamath, and Ayush Sekhari. Hidden poison: Machine unlearning enables camouflaged poisoning attacks. In *NeurIPS ML Safety Workshop*, 2022.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- Cynthia Dwork, Krishnaram Kenthapadi, Frank McSherry, Ilya Mironov, and Moni Naor. Our data, ourselves: Privacy via distributed noise generation. In *Annual international conference on the theory and applications of cryptographic techniques*, pp. 486–503. Springer, 2006.
- Chongyu Fan, Jiancheng Liu, Yihua Zhang, Dennis Wei, Eric Wong, and Sijia Liu. Salun: Empowering machine unlearning via gradient-based weight saliency in both image classification and generation. *arXiv preprint arXiv:2310.12508*, 2023.
- Matt Fredrikson, Somesh Jha, and Thomas Ristenpart. Model inversion attacks that exploit confidence information and basic countermeasures. In *Proceedings of the 22nd ACM SIGSAC conference on computer and communications security*, pp. 1322–1333, 2015.
- Karan Ganju, Qi Wang, Wei Yang, Carl A Gunter, and Nikita Borisov. Property inference attacks on fully connected neural networks using permutation invariant representations. In *Proceedings of the 2018 ACM SIGSAC conference on computer and communications security*, pp. 619–633, 2018.
- Antonio Ginart, Melody Guan, Gregory Valiant, and James Y Zou. Making ai forget you: Data deletion in machine learning. *Advances in neural information processing systems*, 32, 2019.
- Shashwat Goel, Ameya Prabhu, Amartya Sanyal, Ser-Nam Lim, Philip Torr, and Ponnurangam Kumaraguru. Towards adversarial evaluations for inexact machine unlearning. *arXiv preprint arXiv:2201.06640*, 2022.
- 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, 2020.

- Laura Graves, Vineel Nagisetty, and Vijay Ganesh. Amnesiac machine learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 11516–11524, 2021.
- Chuan Guo, Tom Goldstein, Awni Hannun, and Laurens Van Der Maaten. Certified data removal from machine learning models. *arXiv preprint arXiv:1911.03030*, 2019.
- Rob Hall, Alessandro Rinaldo, and Larry Wasserman. Differential privacy for functions and functional data. *arXiv preprint arXiv:1203.2570*, 2012.
- 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*, pp. 770–778, 2016.
- Sorami Hisamoto, Matt Post, and Kevin Duh. Membership inference attacks on sequence-to-sequence models: Is my data in your machine translation system? *Transactions of the Association for Computational Linguistics*, 8:49–63, 2020.
- Nils Homer, Szabolcs Szelinger, Margot Redman, David Duggan, Waibhav Tembe, Jill Muehling, John V Pearson, Dietrich A Stephan, Stanley F Nelson, and David W Craig. Resolving individuals contributing trace amounts of dna to highly complex mixtures using high-density snp genotyping microarrays. *PLoS genetics*, 4(8):e1000167, 2008.
- Chris Jay Hoofnagle, Bart van der Sloot, and Frederik Zuiderveen Borgesius. The european union general data protection regulation: what it is and what it means. *Information & Communications Technology Law*, 28(1):65–98, 2019.
- 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, 2021.
- Saachi Jain, Hadi Salman, Alaa Khaddaj, Eric Wong, Sung Min Park, and Aleksander Madry. A data-based perspective on transfer learning. *arXiv preprint arXiv:2207.05739*, 2022.
- Sue Curry Jansen and Brian Martin. The streisand effect and censorship backfire. *International Journal of Communication*, 9:16, 2015.
- Zhanglong Ji, Zachary C Lipton, and Charles Elkan. Differential privacy and machine learning: a survey and review. *arXiv preprint arXiv:1412.7584*, 2014.
- Jinghan Jia, Jiancheng Liu, Parikshit Ram, Yuguang Yao, Gaowen Liu, Yang Liu, Pranay Sharma, and Sijia Liu. Model sparsification can simplify machine unlearning. *arXiv preprint arXiv:2304.04934*, 2023.
- Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In *International conference on machine learning*, pp. 1885–1894. PMLR, 2017.
- Pang Wei W Koh, Kai-Siang Ang, Hubert Teo, and Percy S Liang. On the accuracy of influence functions for measuring group effects. *Advances in neural information processing systems*, 32, 2019.
- Anastasia Koloskova, Youssef Allouah, Animesh Jha, Rachid Guerraoui, and Sanmi Koyejo. Certified unlearning for neural networks. *arXiv preprint arXiv:2506.06985*, 2025.
- Shuming Kong, Yanyan Shen, and Linpeng Huang. Resolving training biases via influence-based data relabeling. In *International Conference on Learning Representations*, 2021.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Meghdad Kurmanji, Peter Triantafillou, and Eleni Triantafillou. Towards unbounded machine unlearning. *arXiv preprint arXiv:2302.09880*, 2023.

- Ken Lang. Newsweeder: Learning to filter netnews. In *Machine learning proceedings 1995*, pp. 331–339. Elsevier, 1995.
- Klas Leino and Matt Fredrikson. Stolen memories: Leveraging model memorization for calibrated {White-Box} membership inference. In *29th USENIX security symposium (USENIX Security 20)*, pp. 1605–1622, 2020.
- Zheng Li and Yang Zhang. Membership leakage in label-only exposures. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*, pp. 880–895, 2021.
- Yang Liu, Mingyuan Fan, Cen Chen, Ximeng Liu, Zhuo Ma, Li Wang, and Jianfeng Ma. Backdoor defense with machine unlearning. *arXiv preprint arXiv:2201.09538*, 2022a.
- Yi Liu, Lei Xu, Xingliang Yuan, Cong Wang, and Bo Li. The right to be forgotten in federated learning: An efficient realization with rapid retraining. *arXiv preprint arXiv:2203.07320*, 2022b.
- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
- Michał Łukasiak, Srinadh Bhojanapalli, Aditya Menon, and Sanjiv Kumar. Does label smoothing mitigate label noise? In *International Conference on Machine Learning*, pp. 6448–6458. PMLR, 2020.
- Neil G Marchant, Benjamin IP Rubinstein, and Scott Alfeld. Hard to forget: Poisoning attacks on certified machine unlearning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 7691–7700, 2022.
- H Brendan McMahan, Daniel Ramage, Kunal Talwar, and Li Zhang. Learning differentially private recurrent language models. *arXiv preprint arXiv:1710.06963*, 2017.
- Shervin Minaee, Amirali Abdolrashidi, Hang Su, Mohammed Bennamoun, and David Zhang. Biometrics recognition using deep learning: A survey. *Artificial Intelligence Review*, pp. 1–49, 2023.
- Rafael Müller, Simon Kornblith, and Geoffrey E Hinton. When does label smoothing help? *Advances in neural information processing systems*, 32, 2019.
- Milad Nasr, Reza Shokri, and Amir Houmansadr. Comprehensive privacy analysis of deep learning: Passive and active white-box inference attacks against centralized and federated learning. In *2019 IEEE symposium on security and privacy (SP)*, pp. 739–753. IEEE, 2019.
- Seth Neel, Aaron Roth, and Saeed Sharifi-Malvajerdi. Descent-to-delete: Gradient-based methods for machine unlearning. In *Algorithmic Learning Theory*, pp. 931–962. PMLR, 2021.
- 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. *arXiv preprint arXiv:2209.02299*, 2022.
- Stuart L Pardo. The california consumer privacy act: Towards a european-style privacy regime in the united states. *J. Tech. L. & Pol’y*, 23:68, 2018.
- Gabriel Pereyra, George Tucker, Jan Chorowski, Łukasz Kaiser, and Geoffrey Hinton. Regularizing neural networks by penalizing confident output distributions. *arXiv preprint arXiv:1701.06548*, 2017.
- Md Atiqur Rahman, Tanzila Rahman, Robert Laganière, Noman Mohammed, and Yang Wang. Membership inference attack against differentially private deep learning model. *Trans. Data Priv.*, 11(1):61–79, 2018.
- Alexandre Sablayrolles, Matthijs Douze, Cordelia Schmid, Yann Ollivier, and Hervé Jégou. White-box vs black-box: Bayes optimal strategies for membership inference. In *International Conference on Machine Learning*, pp. 5558–5567. PMLR, 2019.

- Prasanna Sattigeri, Soumya Ghosh, Inkit Padhi, Pierre Dognin, and Kush R. Varshney. Fair infinitesimal jackknife: Mitigating the influence of biased training data points without refitting. In *Advances in Neural Information Processing Systems*, 2022.
- Ayush Sekhari, Jayadev Acharya, Gautam Kamath, and Ananda Theertha Suresh. Remember what you want to forget: Algorithms for machine unlearning. *Advances in Neural Information Processing Systems*, 34:18075–18086, 2021.
- Reza Shokri and Vitaly Shmatikov. Privacy-preserving deep learning. In *Proceedings of the 22nd ACM SIGSAC conference on computer and communications security*, pp. 1310–1321, 2015.
- Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models. In *2017 IEEE symposium on security and privacy (SP)*, pp. 3–18. IEEE, 2017.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- Liwei Song and Prateek Mittal. Systematic evaluation of privacy risks of machine learning models. In *30th USENIX security symposium (USENIX security 21)*, pp. 2615–2632, 2021.
- Thomas Steinke, Milad Nasr, and Matthew Jagielski. Privacy auditing with one (1) training run. *Advances in Neural Information Processing Systems*, 36:49268–49280, 2023.
- Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2818–2826, 2016.
- Anvith Thudi, Gabriel Deza, Varun Chandrasekaran, and Nicolas Papernot. Unrolling sgd: Understanding factors influencing machine unlearning. *arXiv preprint arXiv:2109.13398*, 2021.
- Enayat Ullah, Tung Mai, Anup Rao, Ryan A Rossi, and Raman Arora. Machine unlearning via algorithmic stability. In *Conference on Learning Theory*, pp. 4126–4142. PMLR, 2021.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Jialu Wang, Xin Eric Wang, and Yang Liu. Understanding instance-level impact of fairness constraints. In *International Conference on Machine Learning*, pp. 23114–23130. PMLR, 2022a.
- Junxiao Wang, Song Guo, Xin Xie, and Heng Qi. Federated unlearning via class-discriminative pruning. In *Proceedings of the ACM Web Conference 2022*, pp. 622–632, 2022b.
- Alexander Warnecke, Lukas Pirch, Christian Wressnegger, and Konrad Rieck. Machine unlearning of features and labels. *arXiv preprint arXiv:2108.11577*, 2021.
- Jiaheng Wei, Hangyu Liu, Tongliang Liu, Gang Niu, Masashi Sugiyama, and Yang Liu. To smooth or not? when label smoothing meets noisy labels. *Learning*, 1(1):e1, 2021.
- Haotian Ye, Chuanlong Xie, Yue Liu, and Zhenguo Li. Out-of-distribution generalization analysis via influence function. *arXiv preprint arXiv:2101.08521*, 2021.
- Samuel Yeom, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. Privacy risk in machine learning: Analyzing the connection to overfitting. In *2018 IEEE 31st computer security foundations symposium (CSF)*, pp. 268–282. IEEE, 2018.

**Roadmap** The appendix is composed as follows. Section A presents all the notations and their meaning we use in this paper. Section B presents the rest of the Related Work. Section C gives the proof of our theoretical analysis. Section D gives a more detailed explanation of the proposed algorithm. Section E shows the additional experiment results with more details that are not given in the main paper due to the page limit.

## A NOTATION TABLE

The notations used in the paper is summarized in the Table 5.

Table 5: Notation used in this paper

Notations	Description
$K$	The number of classes in the dataset
$\mathcal{D}, \mathcal{X}, \mathcal{Y}$	The general dataset distribution, the feature space and the label space
$D$	The dataset $D \in \mathcal{D}$
$D_{tr}, D_r, D_f$	The training set, remaining set and forgetting set
$\Theta_{\mathcal{M}}$	The distribution of models learned using mechanism $\mathcal{M}$
$\theta$	The model weight
$\theta^*$	The optimal model weight
$\theta_{f,LS}^*$	The optimal model weight trained with $D_f$ whose label is smoothed
$\ \theta\ $	The 2-norm of the model weight
$n$	The size of the dataset
$\varepsilon$	The up-weighted weight of data point $z$ in influence function
$\mathcal{I}(z)$	Influence function of data point $z$
$h_{\theta}$	A function $h$ parameterized by $\theta$
$\ell(h_{\theta}, z_i)$	Loss of $h_{\theta}(x_i)$ and $y_i$
$R_{tr}(\theta)$	The empirical risk of training set when the model weight is $\theta$
$R_f(\theta)$	The empirical risk of forgetting set when the model weight is $\theta$
$R_r(\theta)$	The empirical risk of remaining set when the model weight is $\theta$
$H_{\theta}$	The Hessian matrix w.r.t. $\theta$
$\nabla_{\theta}$	The gradient w.r.t. $\theta$
$B$	Data batch
$B^{LS,\alpha}$	The smoothed batch using $\alpha$
$z_i = (x_i, y_i)$	A data point $z_i$ whose feature is $x_i$ and label is $y_i$
$\mathbf{y}_i$	The one-hot encoded vector form of $y_i$
$\mathbf{y}_i^{GLS,\alpha}$	The smoothed one-hot encoded vector form of $y_i$ where the smooth rate is $\alpha$
$\alpha$	Smooth rate in general label smoothing
$h_{\theta}(x)$	The extracted feature of $x$ from the model parameterized by $\theta$
$\gamma_1, \gamma_2$	The weight of machine learning and machine unlearning on ERM

## B RELATED WORK

**Label Smoothing** (LS) or positive label smoothing (PLS) (Szegedy et al., 2016) is a commonly used regularization method to improve the model performance. Standard training with one-hot labels will lead to overfitting easily. Empirical studies have shown the effectiveness of LS in noisy label (Szegedy et al., 2016; Pereyra et al., 2017; Vaswani et al., 2017; Chorowski & Jaitly, 2016). In addition, LS shows its capability to reduce overfitting, improve generalization, etc. LS can also improve the model calibration (Müller et al., 2019). However, most work on LS is PLS. (Wei et al., 2021) first proposes the concept of negative label smoothing and shows there is a wider feasible domain for the smoothing rate when the rate is negative, expanding the usage of LS.

**Influence Function** is a classic statistical method to track the impact of one training sample. (Koh & Liang, 2017) uses a second-order optimization approximation to evaluate the impact of a training sample. Additionally, it can also be used to identify the importance of the training groups (Basu et al., 2020; Koh et al., 2019). The influence function is widely used in many machine-learning

tasks, such as data bias solution (Brunet et al., 2019; Kong et al., 2021), fairness (Sattigeri et al., 2022; Wang et al., 2022a), security (Liu et al., 2022a), transfer learning (Jain et al., 2022), out-of-distribution generalization (Ye et al., 2021), etc. The approach also plays an important role as the algorithm backbone in the MU tasks (Jia et al., 2023; Warnecke et al., 2021; Izzo et al., 2021).

**Differential Privacy** (DP) is a mathematical framework designed to quantify and mitigate privacy risks in machine learning models. It ensures that the inclusion or exclusion of a single data point in a dataset does not significantly affect the model’s output, thus protecting individual data points from being inferred by adversaries (Dwork et al., 2006). In machine learning, DP mechanisms such as noise addition and gradient clipping are employed during the training process to provide formal privacy guarantees while maintaining model utility (Abadi et al., 2016). These techniques help balance the trade-off between data privacy and model performance, making DP a cornerstone of privacy-preserving machine learning (Shokri & Shmatikov, 2015; McMahan et al., 2017).

A multitude of **privacy risk assessment** tools have been proposed to gauge the degree of leakage associated with the training data. Specifically targeted at the training data, model attacks are often used as a proxy metric for privacy leakage in pretrained models. For example, model inversion attacks are designed to extract aggregate information about specific sub-classes rather than individual samples (Fredrikson et al., 2015). Data extraction attacks aim to reverse engineer individual samples used during training (Carlini et al., 2021), while property inference attacks focus on inferring properties of the training data (Ganju et al., 2018).

More relevant to the current work are **Membership Inference Attacks** (MIA), which predict whether a particular sample was used to train the model. First introduced by (Homer et al., 2008), membership attack algorithms were later formalized in the context of DP, enabling privacy attacks and defenses for machine learning models (Rahman et al., 2018). (Shokri et al., 2017) introduced MIA based on the assumption of adversarial queries to the target model. By training a reference attack model (shadow model) based on the model inference response, this type of MIA has proven to be powerful in scenarios such as white-box (Leino & Fredrikson, 2020; Nasr et al., 2019; Sablayrolles et al., 2019), black-box (Chen et al., 2021a; Hisamoto et al., 2020; Song & Mittal, 2021), and label-only (Choquette-Choo et al., 2021; Li & Zhang, 2021) access. However, most MIA mechanisms often require training a large number of shadow models with diverse subsets of queries, making them prohibitively expensive. As a result, some recent works have focused on developing cheaper MIA mechanisms (Steinke et al., 2023).

**Basics of Influence Function** Given a dataset  $D = \{z_i : (x_i, y_i)\}_{i=1}^n$  and a function  $h$  parameterized by  $\theta$  which maps from the input feature space  $\mathcal{X}$  to the output space  $\mathcal{Y}$ . Recall the standard empirical risk minimization is written as:

$$\theta^* = \arg \min_{\theta} \frac{1}{n} \sum_{z \in D} \ell(h_{\theta}, z). \quad (9)$$

To find the impact of a training point  $\hat{z}$ , we up-weight its weight by an infinitesimal amount  $\varepsilon^1$ . The new model parameter  $\theta_{\{\hat{z}\}}^{\varepsilon}$  can be obtained from

$$\theta_{\{\hat{z}\}}^{\varepsilon} = \arg \min_{\theta} \frac{1}{n} \sum_{z \in D} \ell(h_{\theta}, z) + \varepsilon \cdot \ell(h_{\theta}, \hat{z}) \quad (10)$$

When  $\varepsilon = -\frac{1}{n}$ , it is indicating removing  $\hat{z}$ . According to (Koh & Liang, 2017),  $\theta_{\{\hat{z}\}}^{\varepsilon}$  can be approximated by using the first-order Taylor series expansion as

$$\theta_{\{\hat{z}\}}^{\varepsilon} \approx \theta^* - \varepsilon \cdot H_{\theta^*}^{-1} \cdot \nabla_{\theta} \ell(h_{\theta^*}, \hat{z}), \quad (11)$$

where  $H_{\theta^*}$  is the Hessian with respect to (w.r.t.)  $\theta^*$ . The change of  $\theta$  due to changing the weight can be given using the influence function  $\mathcal{I}(\hat{z})$  as

$$\Delta \theta = \theta_{\{\hat{z}\}}^{\varepsilon} - \theta^* = \mathcal{I}(\hat{z}) = \left. \frac{d\theta_{\{\hat{z}\}}^{\varepsilon}}{d\varepsilon} \right|_{\varepsilon=0} = -H_{\theta^*}^{-1} \cdot \nabla_{\theta} \ell(h_{\theta^*}, \hat{z}).$$

<sup>1</sup>To distinguish it from the  $\epsilon$  in differential privacy, we use  $\varepsilon$  here.

## C PROOFS

### C.1 PROOF FOR THEOREM 1

*Proof.* For  $p(x)$ , the Taylor expansion at  $x = a$  is

$$p(x) = p(a) + \frac{p'(a)}{1}(x - a) + o \quad (12)$$

Here,  $p(\boldsymbol{\theta}) = \nabla R_{tr}(\boldsymbol{\theta}) + \varepsilon \sum_{D_f} \nabla \ell(h_{\boldsymbol{\theta}}, z_i^f)$  so we have

$$p(\boldsymbol{\theta}) = \nabla R_{tr}(a) + \varepsilon \sum_{z^f \in D_f} \nabla \ell(h_a, z^f) + \left[ \nabla^2 R_{tr}(a) + \varepsilon \sum_{z^f \in D_f} \nabla^2 \ell(h_a, z^f) \right] (\boldsymbol{\theta} - a) + o \quad (13)$$

For Eq. (2), we expand  $p(\boldsymbol{\theta}_f^*)$  at  $\boldsymbol{\theta} = \boldsymbol{\theta}_{tr}^*$  as

$$\begin{aligned} p(\boldsymbol{\theta}_f^*) &= \nabla R_{tr}(\boldsymbol{\theta}_{tr}^*) + \varepsilon \sum_{z^f \in D_f} \nabla \ell(h_{\boldsymbol{\theta}_{tr}^*}, z^f) \\ &+ \left[ \nabla^2 R_{tr}(\boldsymbol{\theta}_{tr}^*) + \varepsilon \sum_{z^f \in D_f} \nabla^2 \ell(h_{\boldsymbol{\theta}_{tr}^*}, z^f) \right] (\boldsymbol{\theta}_f^* - \boldsymbol{\theta}_{tr}^*) + o = 0 \end{aligned} \quad (14)$$

Since we have  $\nabla R_{tr}(\boldsymbol{\theta}_{tr}^*) = 0$  and ignore  $o$ , we can get the approximation as

$$\boldsymbol{\theta}_f^* - \boldsymbol{\theta}_{tr}^* \approx - \left[ \sum_{z^{tr} \in D_{tr}} \nabla^2 \ell(h_{\boldsymbol{\theta}_{tr}^*}, z^{tr}) + \varepsilon \sum_{z^f \in D_f} \nabla^2 \ell(h_{\boldsymbol{\theta}_{tr}^*}, z^f) \right]^{-1} \left[ \varepsilon \sum_{z^f \in D_f} \nabla \ell(h_{\boldsymbol{\theta}_{tr}^*}, z^f) \right] \quad (15)$$

Similarly, we can expand  $q(\boldsymbol{\theta}_{tr}^*) = \nabla R_{tr}(\boldsymbol{\theta}_{tr}^*)$  at  $\boldsymbol{\theta} = \boldsymbol{\theta}_r^*$  as

$$\begin{aligned} q(\boldsymbol{\theta}_{tr}^*) &= \sum_{z^{tr} \in D_{tr}} \nabla \ell(h_{\boldsymbol{\theta}_r^*}, z^{tr}) + \sum_{z^{tr} \in D_{tr}} \nabla^2 \ell(h_{\boldsymbol{\theta}_r^*}, z^{tr}) (\boldsymbol{\theta}_{tr}^* - \boldsymbol{\theta}_r^*) \approx 0 \\ \boldsymbol{\theta}_r^* - \boldsymbol{\theta}_{tr}^* &\approx \left[ \sum_{z^{tr} \in D_{tr}} \nabla^2 \ell(h_{\boldsymbol{\theta}_r^*}, z^{tr}) \right]^{-1} \sum_{z^{tr} \in D_{tr}} \nabla \ell(h_{\boldsymbol{\theta}_r^*}, z^{tr}) \end{aligned} \quad (16)$$

Because of gradient ascent,  $\varepsilon = -1$  and we have

$$\begin{aligned} \boldsymbol{\theta}_r^* - \boldsymbol{\theta}_f^* &= \boldsymbol{\theta}_r^* - \boldsymbol{\theta}_{tr}^* - (\boldsymbol{\theta}_{tr}^* - \boldsymbol{\theta}_f^*) = \underbrace{\left( \sum_{z^{tr} \in D_{tr}} \nabla^2 \ell(h_{\boldsymbol{\theta}_r^*}, z^{tr}) \right)^{-1} \sum_{z^{tr} \in D_{tr}} \nabla \ell(h_{\boldsymbol{\theta}_r^*}, z^{tr})}_{\Delta \boldsymbol{\theta}_r} \\ &- \underbrace{\left( \sum_{z^r \in D_r} \nabla^2 \ell(h_{\boldsymbol{\theta}_{tr}^*}, z^r) \right)^{-1} \sum_{z^f \in D_f} \nabla \ell(h_{\boldsymbol{\theta}_{tr}^*}, z^f)}_{\Delta \boldsymbol{\theta}_f} \end{aligned} \quad (17)$$

Thus,  $\|\boldsymbol{\theta}_r^* - \boldsymbol{\theta}_f^*\| = 0$  if and only if  $\Delta \boldsymbol{\theta}_f = \Delta \boldsymbol{\theta}_r$ , where

$$\sum_{z^{tr} \in D_{tr}} \nabla \ell(h_{\boldsymbol{\theta}_r^*}, z^{tr}) = \underbrace{\left[ \sum_{z^{tr} \in D_{tr}} \nabla^2 \ell(h_{\boldsymbol{\theta}_r^*}, z^{tr}) \right] \left[ \sum_{z^r \in D_r} \nabla^2 \ell(h_{\boldsymbol{\theta}_{tr}^*}, z^r) \right]^{-1}}_{H(\boldsymbol{\theta}_r^*, \boldsymbol{\theta}_{tr}^*)} \sum_{z^f \in D_f} \nabla \ell(h_{\boldsymbol{\theta}_{tr}^*}, z^f) \quad (18)$$

□

## C.2 ERROR ANALYSIS IN THEOREM 1

If we do not ignore the Lagrange remainder in Eq. 14 and 16 and denote them as  $e_r$  and  $e_f$ , Eq. 14 and 16 become

$$\begin{aligned} p(\boldsymbol{\theta}_f^*) &= \nabla R_{tr}(\boldsymbol{\theta}_{tr}^*) + \varepsilon \sum_{z^f \in D_f} \nabla \ell(h_{\boldsymbol{\theta}_{tr}^*}, z^f) \\ &+ \left[ \nabla^2 R_{tr}(\boldsymbol{\theta}_{tr}^*) + \varepsilon \sum_{z^f \in D_f} \nabla^2 \ell(h_{\boldsymbol{\theta}_{tr}^*}, z^f) \right] (\boldsymbol{\theta}_f^* - \boldsymbol{\theta}_{tr}^*) + e_r = 0 \end{aligned} \quad (19)$$

$$q(\boldsymbol{\theta}_{tr}^*) = \sum_{z^{tr} \in D_{tr}} \nabla \ell(h_{\boldsymbol{\theta}_r^*}, z^{tr}) + \sum_{z^{tr} \in D_{tr}} \nabla^2 \ell(h_{\boldsymbol{\theta}_r^*}, z^{tr}) (\boldsymbol{\theta}_{tr}^* - \boldsymbol{\theta}_r^*) + e_f = 0 \quad (20)$$

, respectively. Thus,

$$\boldsymbol{\theta}_r^* - \boldsymbol{\theta}_f^* = (\boldsymbol{\theta}_r^* - \boldsymbol{\theta}_{tr}^*) - (\boldsymbol{\theta}_f^* - \boldsymbol{\theta}_{tr}^*) \quad (21)$$

$$= (\Delta \boldsymbol{\theta}_r + e_r) - (\Delta \boldsymbol{\theta}_f + e_f) = (\Delta \boldsymbol{\theta}_r - \Delta \boldsymbol{\theta}_f) + (e_r - e_f). \quad (22)$$

We now bound the error of using the linearized difference  $\Delta \boldsymbol{\theta}_r - \Delta \boldsymbol{\theta}_f$  to approximate  $\boldsymbol{\theta}_r^* - \boldsymbol{\theta}_f^*$ .

$$\boldsymbol{\theta}_r^* - \boldsymbol{\theta}_f^* - (\Delta \boldsymbol{\theta}_r - \Delta \boldsymbol{\theta}_f) = e_r - e_f, \quad (23)$$

and hence

$$\|\boldsymbol{\theta}_r^* - \boldsymbol{\theta}_f^* - (\Delta \boldsymbol{\theta}_r - \Delta \boldsymbol{\theta}_f)\| = \|e_r - e_f\| \leq \|e_r\| + \|e_f\|. \quad (24)$$

Assume that  $q(\boldsymbol{\theta}) = \nabla R_{tr}(\boldsymbol{\theta})$  and  $p(\boldsymbol{\theta}) = \nabla R_{tr}(\boldsymbol{\theta}) - \nabla R_f(\boldsymbol{\theta})$  have Lipschitz-continuous Hessians with constants  $L_q$  and  $L_p$ , respectively, i.e.,

$$\|\nabla^2 q(\boldsymbol{\theta}_1) - \nabla^2 q(\boldsymbol{\theta}_2)\| \leq L_q \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|, \quad (25)$$

$$\|\nabla^2 p(\boldsymbol{\theta}_1) - \nabla^2 p(\boldsymbol{\theta}_2)\| \leq L_p \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|. \quad (26)$$

Then standard Taylor bounds imply

$$\|r_q\| \leq \frac{L_q}{2} \|\boldsymbol{\theta}_{tr}^* - \boldsymbol{\theta}_r^*\|^2, \quad (27)$$

$$\|r_p\| \leq \frac{L_p}{2} \|\boldsymbol{\theta}_f^* - \boldsymbol{\theta}_{tr}^*\|^2. \quad (28)$$

Using  $e_r = -H_r^{-1} r_q$  and  $e_f = -H_f^{-1} r_p$ , we obtain

$$\|e_r\| \leq \|H_r^{-1}\| \|r_q\| \leq \frac{L_q}{2} \|H_r^{-1}\| \|\boldsymbol{\theta}_r^* - \boldsymbol{\theta}_{tr}^*\|^2, \quad (29)$$

$$\|e_f\| \leq \|H_f^{-1}\| \|r_p\| \leq \frac{L_p}{2} \|H_f^{-1}\| \|\boldsymbol{\theta}_f^* - \boldsymbol{\theta}_{tr}^*\|^2. \quad (30)$$

Therefore, the approximation error satisfies

$$\|\boldsymbol{\theta}_r^* - \boldsymbol{\theta}_f^* - (\Delta \boldsymbol{\theta}_r - \Delta \boldsymbol{\theta}_f)\| \leq \frac{L_q}{2} \|H_r^{-1}\| \|\boldsymbol{\theta}_r^* - \boldsymbol{\theta}_{tr}^*\|^2 + \frac{L_p}{2} \|H_f^{-1}\| \|\boldsymbol{\theta}_f^* - \boldsymbol{\theta}_{tr}^*\|^2. \quad (31)$$

## C.3 PROOF FOR THEOREM 2

*Proof.* Recall the loss calculation in label smoothing and we have

$$\ell(h_{\boldsymbol{\theta}}, z^{\text{GLS}, \alpha}) = \left(1 + \frac{1-K}{K} \alpha\right) \ell(h_{\boldsymbol{\theta}}, (x, y)) + \frac{\alpha}{K} \sum_{y' \in \mathcal{Y} \setminus y} \ell(h_{\boldsymbol{\theta}}, (x, y')), \quad (32)$$

where we use notations  $\ell(h_{\boldsymbol{\theta}}, (x, y)) := \ell(h_{\boldsymbol{\theta}}, z)$  to specify the loss of an example  $z = \{x, y\}$  existing in the dataset and  $\ell(h_{\boldsymbol{\theta}}, (x, y'))$  to denote the loss of an example when its label is replaced

with  $y'$ .  $\nabla_{\theta}\ell(h_{\theta}, (x, y))$  is the gradient of the target label and  $\sum_{y' \in \mathcal{Y} \setminus y} \nabla_{\theta}\ell(h_{\theta}, (x, y'))$  is the sum of the gradient of non-target labels.

With label smoothing in Eq. (32), Eq. (17) becomes

$$\begin{aligned} \theta_r^* - \theta_{f,LS}^* &\approx \Delta\theta_r + \left(1 + \frac{1-K}{K}\alpha\right) \cdot (-\Delta\theta_f) + \frac{1-K}{K}\alpha \cdot \Delta\theta_n \\ &= \Delta\theta_r - \Delta\theta_f + \frac{1-K}{K}\alpha \cdot (\Delta\theta_n - \Delta\theta_f) \end{aligned} \quad (33)$$

where

$$\begin{aligned} \Delta\theta_r &:= \left[ \sum_{z^{tr} \in D_{tr}} \nabla_{\theta}^2 \ell(h_{\theta_r^*}, z^{tr}) \right]^{-1} \sum_{z^{tr} \in D_{tr}} \nabla_{\theta} \ell(h_{\theta_r^*}, z^{tr}) \\ \Delta\theta_f &:= \left[ \sum_{z^r \in D_r} \nabla_{\theta}^2 \ell(h_{\theta_r^*}, z^r) \right]^{-1} \sum_{z^f \in D_f} \nabla_{\theta} \ell(h_{\theta_r^*}, z^f) \end{aligned}$$

as given in Eq. (17). So we have

$$\theta_r^* - \theta_{f,LS}^* \approx \Delta\theta_r - \Delta\theta_f + \frac{1-K}{K}\alpha \cdot (\Delta\theta_n - \Delta\theta_f) \quad (34)$$

where

$$\Delta\theta_n := \frac{1}{K-1} \left[ \sum_{z^r \in D_r} \nabla_{\theta}^2 \ell(h_{\theta_r^*}, z^r) \right]^{-1} \sum_{z^f \in D_f} \nabla_{\theta} \sum_{y' \in \mathcal{Y} \setminus y^f} \ell(h_{\theta_r^*}, (x^f, y'))$$

When we have

$$\langle \Delta\theta_r - \Delta\theta_f, \Delta\theta_n - \Delta\theta_f \rangle \leq 0, \quad (35)$$

$\alpha < 0$  can help with MU, making

$$\|\theta_r^* - \theta_{f,NLS}^*\| \leq \|\theta_r^* - \theta_f^*\| \quad (36)$$

□

#### C.4 PROOF FOR THEOREM 3

*Proof.* When the optimization is gradient ascent (GA) with negative label smoothing (NLS), Eq. (6) can be written as

$$\ell(h_{\theta}, z^{\text{NLS},\alpha}) = - \left(1 + \frac{1-K}{K}\alpha\right) \cdot \ell(h_{\theta}, (x, y)) - \frac{\alpha}{K} \sum_{y' \in \mathcal{Y} \setminus y} \ell(h_{\theta}, (x, y')), \alpha < 0, \quad (37)$$

Recall  $R_{tr}(\theta) = \sum_{z^{tr} \in D_{tr}} \ell(h_{\theta}, z^{tr})$ . Denote by  $R_f^{\text{NLS}}(\theta; \alpha) = \sum_{z^f \in D_f} \ell(h_{\theta}, z^f)$ ,  $\alpha < 0$  the empirical risk of forgetting data with NLS. After MU with label smoothing on  $D_f$  by gradient ascent, the resulting model can be seen as minimizing the risk  $\gamma_1 \cdot R_{tr}(\theta) - \gamma_2 \cdot R_f^{\text{NLS}}(\theta; \alpha)$ , which is a weighted combination of the risk from two phases: 1) machine learning on  $D_{tr}$  with weight  $\gamma_1 > 0$  and 2) machine unlearning on  $D_f$  with weight  $\gamma_2 > 0$ . Consider an example  $(x, y)$  in the forgetting dataset. The loss of this example is:

$$\begin{aligned} \gamma_1 \ell(h_{\theta}, (x, y)) - \gamma_2 \ell(h_{\theta}, z^{\text{GLS},\alpha}) &= \left[ \gamma_1 - \gamma_2 \left(1 + \frac{1-K}{K}\alpha\right) \right] \cdot \ell(h_{\theta}, (x, y)) \\ &\quad - \frac{\alpha}{K} \gamma_2 \sum_{y' \in \mathcal{Y} \setminus y} \ell(h_{\theta}, (x, y')). \end{aligned}$$

When  $[\gamma_1 - \gamma_2 (1 + \frac{1-K}{K}\alpha)] > 0$ , the optimal solution by minimizing this loss is

$$\mathbb{P}(\mathcal{M}(y) = y^{\text{pred}}) = \begin{cases} \frac{\gamma_1 - \gamma_2 (1 + \frac{1-K}{K}\alpha)}{(\gamma_1 - \gamma_2 (1 + \frac{1-K}{K}\alpha)) - \frac{K-1}{K}\alpha\gamma_2}, & \text{if } y^{\text{pred}} = y, \\ \frac{-\frac{K}{K}\gamma_2}{(\gamma_1 - \gamma_2 (1 + \frac{1-K}{K}\alpha)) - \frac{K-1}{K}\alpha\gamma_2}, & \text{if } y^{\text{pred}} \neq y. \end{cases}$$

Accordingly, for another label  $y'$ , we have

$$\mathbb{P}(\mathcal{M}(y') = y^{\text{pred}}) = \begin{cases} \frac{\gamma_1 - \gamma_2 (1 + \frac{1-K}{K}\alpha)}{(\gamma_1 - \gamma_2 (1 + \frac{1-K}{K}\alpha)) - \frac{K-1}{K}\alpha\gamma_2}, & \text{if } y^{\text{pred}} = y', \\ \frac{-\frac{K}{K}\gamma_2}{(\gamma_1 - \gamma_2 (1 + \frac{1-K}{K}\alpha)) - \frac{K-1}{K}\alpha\gamma_2}, & \text{if } y^{\text{pred}} \neq y'. \end{cases}$$

Then the quotient of two probabilities can be upper bounded by:

$$\log\left(\frac{\mathbb{P}(\mathcal{M}(y) = y^{\text{pred}})}{\mathbb{P}(\mathcal{M}(y') = y^{\text{pred}})}\right) \leq \left| \log\left(\frac{\gamma_1 - \gamma_2 (1 + \frac{1-K}{K}\alpha)}{-\frac{\alpha}{K} \cdot \gamma_2}\right) \right| = \left| \log\left(\frac{K}{\alpha} (1 - \frac{\gamma_1}{\gamma_2}) + 1 - K\right) \right| = \epsilon.$$

□

## D THE DETAILS OF ALGORITHM

### D.1 ALGORITHM DETAILS

We provide a more detailed explanation of UGradSL and UGradSL+ in Algorithm 1 here. For UGradSL+, we first sample a batch  $B_r = \{z_i^r : (x_i^r, y_i^r)\}_{i=1}^{n_{B_r}}$  from  $D_r$  (Line 3-4). Additionally, we sample a batch  $B_f = \{z_i^f : (x_i^f, y_i^f)\}_{i=1}^{n_{B_f}}$  from  $D_f$  where  $n_{B_r} = n_{B_f}$  (Line 5). We compute the distance  $d(z_i^r, z_i^f) \in [0, 1]$  for each  $(z_i^r, z_i^f)$  pair where  $z_i^r \in B_r$  and  $z_i^f \in B_f$  (Line 6). For each  $z_i^f$ , we count the number of  $z_i^r$  whose  $d(z_i^r, z_i^f) < \beta$ , where  $\beta$  is the distance threshold. This count is denoted by  $c_i^f$  (Line 7). Then we get the smooth rate by normalizing the count as  $\alpha_i = c_i^f / |B_f|$ , where  $\alpha_i \in [0, 1]$  (Line 8). GA with NLS is to decrease the model confidence of  $D_f$ . The larger the absolute value of  $\alpha_i$ , the lower confidence will be given. Our intuition is that a smaller  $d(z_i^r, z_i^f)$  means  $z_i^r$  is more similar to  $D_r$  and the confidence of  $z_i^f$  should not be decreased too much. The distances we use is cosine distance. UGradSL is similar and the difference is the dataset replacement. For each epoch, UGradSL+ is terminated after completing the iterations on  $D_r$ , while UGradSL is terminated after completing the iterations on  $D_f$ .

### D.2 ALGORITHM EXPLANATION

In the self-adaptive version of UGradSL+, the label smoothing rate for each forgetting sample is computed dynamically from its proximity to the retained data in feature space. For each iteration, the algorithm samples a batch of retained examples  $B_r$  and a batch of forgetting examples  $B_f$  with equal size, extracts their features  $\{z_i^r\}$  and  $\{z_j^f\}$ , and computes the **feature distance**  $d(z_i^r, z_j^f)$  for every retained-forgetting pair. Then, for each forgetting feature  $z_j^f$ , it counts how many retained features fall within a distance threshold  $\beta$ , denoted as  $c_j^f$ . This count is normalized by the batch size  $|B_f|$  to obtain the adaptive smoothing rate  $\alpha_j = c_j^f / |B_f|$ . As a result, forgetting samples that are close to many retained samples (i.e., highly entangled in representation space) receive a higher smoothing rate and are updated more conservatively, while those that are far from retained data get a lower smoothing rate (possibly zero) and can be pushed away more aggressively during unlearning.

### D.3 ADDITIONAL RESULTS

As mentioned in Section 4, to avoid the smooth rate selection, we propose a self-adaptive smooth rate version. We compare the performance with and without self-adaptive smooth on CIFAR-10 and SVHN. The forgetting scenario is random forgetting. The results are given in Table 10.

#### D.4 COMPLEXITY ANALYSIS

Compared with the fixed  $\alpha$ , the additional computation from the adaptive version is the distance calculation. The code we compute the distance is given below. All computations are implemented as **batched GPU tensor operations without any explicit Python loops**. We assume the feature from  $D_r$  and  $D_f$  are both in  $\mathbb{R}^{n \times d}$ , where  $n$  is the batch size and  $d$  is the feature dimension.

For the FLOP count,

- The two normalization operations cost approximately  $6nd$  FLOPs in total, since normalizing a single  $n \times d$  tensor requires about  $3nd$  FLOPs (square, sum, and division).
- Computing the cosine similarity matrix costs about  $2n^2d$  FLOPs, as each of the  $n^2$  entries is a dot product between two  $d$ -dimensional vectors.
- Converting similarity to distance and applying the threshold require about  $2n^2$  and  $n^2$  FLOPs, respectively.
- The density computation costs about  $n^2$  FLOPs for forming the mask and  $n$  FLOPs for the length normalization.

Overall, the total FLOP count is  $6nd + 2n^2d + 4n^2 + n$ , which is dominated by the  $O(n^2d)$  cosine-similarity term. For our typical setting  $n = 64$  and  $d = 512$ , this corresponds to roughly  $4.4 \times 10^6$  FLOPs. Compared with the FP32 peak throughput of an A6000 GPU (38.71 TFLOPS), this overhead is negligible relative to the usual forward/backward passes.

For memory usage, the additional GPU tensors have the following shapes:

- Each features:  $n \times d$
- Each normalized features:  $n \times d$
- The cosine similarity, cosine distance and the filtered mask:  $n \times n$
- The density:  $n$

Assuming FP32 (4 bytes) for all tensors, the peak extra memory is at most  $4(4nd + 3n^2 + n)$  bytes, which is 561,408 bytes ( $\approx 0.5$  MiB) for  $n = 64$  and  $d = 512$ . This is negligible compared with the model parameters, so the memory overhead can also be safely ignored.

```

1 forget_norm = F.normalize(forget_feature, p=2, dim=1)
2 retain_norm = F.normalize(retain_feature, p=2, dim=1)
3 # cosine similarity (batch x batch)
4 cos_sim = forget_norm @ retain_norm.T
5 # convert to distance in [0,1]
6 cos_dist = (1 - cos_sim) / 2 # shape: [batch, batch]
7 # --- threshold and count ---
8 threshold = 0.2 # example threshold in [0,1]
9 # boolean matrix: True = close
10 close_mask = cos_dist < threshold # [batch, batch]
11 density = close_mask.sum(dim=0) / len(forget_feature)

```

#### D.5 ABLATION STUDY

By default, we adopt cosine distance because it naturally lies in  $[0, 1]$ , and we set the threshold  $\beta$  to the median of all pairwise distances. We conduct an ablation study on different distance metrics and thresholds for random forgetting on CIFAR-10, where the forgetting set size is 10% of the training data. For comparison, we also evaluate Euclidean distance. The results for cosine and Euclidean

Table 6: The ablation studies of threshold  $\beta$  and different distance functions of UGradSL for the random forgetting on CIFAR-10 and the size of forgetting set is 10% of the training set. The first row is the results for retraining for reference.

$\beta$	Distance	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap ( $\downarrow$ )
-	-	8.07	17.41	100.00	91.61	-
Median	Cosine	6.04 $\pm$ 0.11	13.75 $\pm$ 0.32	99.11 $\pm$ 0.01	92.07 $\pm$ 0.02	1.76
	Euclidean	8.59 $\pm$ 1.85	17.30 $\pm$ 0.98	94.39 $\pm$ 1.15	88.97 $\pm$ 1.12	2.22
0.1	Cosine	6.50 $\pm$ 0.14	14.76 $\pm$ 1.52	95.64 $\pm$ 0.23	89.91 $\pm$ 0.17	2.57
	Euclidean	6.68 $\pm$ 0.88	14.69 $\pm$ 1.66	95.34 $\pm$ 0.79	89.90 $\pm$ 0.69	2.62
0.2	Cosine	7.01 $\pm$ 0.67	15.86 $\pm$ 0.86	95.18 $\pm$ 0.44	89.69 $\pm$ 0.19	2.34
	Euclidean	6.82 $\pm$ 0.44	15.81 $\pm$ 0.70	95.58 $\pm$ 0.73	90.02 $\pm$ 0.57	2.21
0.3	Cosine	7.01 $\pm$ 0.98	15.13 $\pm$ 1.26	95.24 $\pm$ 0.99	89.76 $\pm$ 0.41	2.49
	Euclidean	7.32 $\pm$ 1.06	16.45 $\pm$ 2.08	94.68 $\pm$ 0.89	89.16 $\pm$ 0.33	2.37
0.4	Cosine	7.91 $\pm$ 0.26	15.69 $\pm$ 1.11	94.69 $\pm$ 0.51	89.07 $\pm$ 0.29	2.43
	Euclidean	6.24 $\pm$ 0.21	14.16 $\pm$ 0.12	95.75 $\pm$ 0.40	90.13 $\pm$ 0.13	2.70
0.5	Cosine	7.61 $\pm$ 0.66	16.50 $\pm$ 1.68	95.03 $\pm$ 0.36	89.69 $\pm$ 0.72	2.07
	Euclidean	8.27 $\pm$ 1.33	16.44 $\pm$ 1.83	94.67 $\pm$ 1.33	89.03 $\pm$ 1.28	2.27
0.6	Cosine	8.76 $\pm$ 0.28	16.53 $\pm$ 1.88	94.31 $\pm$ 0.61	88.54 $\pm$ 0.50	2.58
	Euclidean	8.67 $\pm$ 0.28	17.01 $\pm$ 2.43	94.34 $\pm$ 0.16	88.93 $\pm$ 0.30	2.34
0.7	Cosine	9.88 $\pm$ 1.05	18.33 $\pm$ 2.82	93.55 $\pm$ 0.92	88.08 $\pm$ 0.42	3.18
	Euclidean	9.61 $\pm$ 0.86	17.93 $\pm$ 2.33	94.11 $\pm$ 0.49	88.69 $\pm$ 0.19	2.72
0.8	Cosine	9.61 $\pm$ 1.12	16.91 $\pm$ 1.51	93.68 $\pm$ 1.20	88.48 $\pm$ 0.76	2.87
	Euclidean	9.75 $\pm$ 0.17	16.79 $\pm$ 0.52	93.87 $\pm$ 0.02	88.34 $\pm$ 0.39	2.93
0.9	Cosine	9.19 $\pm$ 0.66	17.84 $\pm$ 0.72	94.19 $\pm$ 0.50	88.51 $\pm$ 0.84	2.63
	Euclidean	9.76 $\pm$ 0.49	18.61 $\pm$ 0.65	93.90 $\pm$ 0.39	88.47 $\pm$ 0.35	3.03
1.0	Cosine	9.39 $\pm$ 0.07	16.94 $\pm$ 0.26	94.26 $\pm$ 0.33	88.74 $\pm$ 0.22	2.60
	Euclidean	10.41 $\pm$ 0.24	19.16 $\pm$ 1.08	93.50 $\pm$ 0.63	88.21 $\pm$ 0.34	3.50

Table 7: The ablation studies of threshold  $\beta$  and different distance functions of UGradSL+ for the random forgetting on CIFAR-10 and the size of forgetting set is 10% of the training set. The first row is the results for retraining for reference.

$\beta$	Distance	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap ( $\downarrow$ )
-	-	8.07	17.41	100.00	91.61	-
Median	Cosine	7.54 $\pm$ 0.43	13.57 $\pm$ 0.12	99.67 $\pm$ 0.00	92.97 $\pm$ 0.17	1.52
	Euclidean	11.21 $\pm$ 0.21	21.02 $\pm$ 2.23	94.35 $\pm$ 0.22	88.58 $\pm$ 0.26	3.86
0.1	Cosine	7.79 $\pm$ 0.52	17.04 $\pm$ 0.61	95.84 $\pm$ 0.27	90.10 $\pm$ 0.47	1.58
	Euclidean	7.30 $\pm$ 0.62	16.42 $\pm$ 0.66	96.16 $\pm$ 0.94	90.46 $\pm$ 0.91	1.69
0.2	Cosine	8.38 $\pm$ 0.19	17.46 $\pm$ 1.09	95.38 $\pm$ 0.34	89.56 $\pm$ 0.53	1.76
	Euclidean	7.80 $\pm$ 0.76	16.55 $\pm$ 1.91	95.75 $\pm$ 1.04	89.80 $\pm$ 0.50	1.80
0.3	Cosine	8.27 $\pm$ 0.65	18.19 $\pm$ 0.29	95.94 $\pm$ 0.84	90.18 $\pm$ 0.62	1.62
	Euclidean	7.68 $\pm$ 0.65	17.28 $\pm$ 0.52	95.85 $\pm$ 0.75	90.25 $\pm$ 0.55	1.51
0.4	Cosine	8.49 $\pm$ 0.28	17.92 $\pm$ 0.52	95.85 $\pm$ 0.20	90.09 $\pm$ 0.03	1.66
	Euclidean	8.38 $\pm$ 0.60	17.86 $\pm$ 0.89	95.60 $\pm$ 0.78	90.06 $\pm$ 0.57	1.68
0.5	Cosine	9.23 $\pm$ 0.89	16.81 $\pm$ 1.66	95.46 $\pm$ 0.62	89.79 $\pm$ 0.86	2.03
	Euclidean	8.98 $\pm$ 0.69	16.77 $\pm$ 1.62	95.39 $\pm$ 1.01	89.34 $\pm$ 1.17	2.11
0.6	Cosine	9.95 $\pm$ 0.64	19.90 $\pm$ 0.95	95.47 $\pm$ 0.12	89.82 $\pm$ 0.30	2.67
	Euclidean	10.00 $\pm$ 0.10	19.00 $\pm$ 1.92	95.15 $\pm$ 0.26	89.53 $\pm$ 0.28	2.61
0.7	Cosine	11.81 $\pm$ 0.74	20.67 $\pm$ 2.62	94.25 $\pm$ 0.76	88.78 $\pm$ 1.02	3.90
	Euclidean	11.25 $\pm$ 0.59	21.54 $\pm$ 1.12	94.69 $\pm$ 0.71	89.05 $\pm$ 0.71	3.79
0.8	Cosine	13.06 $\pm$ 0.53	18.81 $\pm$ 0.81	92.89 $\pm$ 0.69	87.29 $\pm$ 0.75	4.45
	Euclidean	12.07 $\pm$ 0.45	19.23 $\pm$ 2.00	93.81 $\pm$ 0.95	88.34 $\pm$ 1.00	3.82
0.9	Cosine	11.75 $\pm$ 0.09	21.02 $\pm$ 1.43	94.34 $\pm$ 0.38	88.81 $\pm$ 0.31	3.94
	Euclidean	12.01 $\pm$ 1.12	21.49 $\pm$ 1.17	94.26 $\pm$ 1.08	88.74 $\pm$ 0.88	4.16
1.0	Cosine	11.48 $\pm$ 0.06	20.59 $\pm$ 2.63	94.19 $\pm$ 0.56	88.82 $\pm$ 0.38	3.80
	Euclidean	11.79 $\pm$ 0.37	17.35 $\pm$ 0.85	94.37 $\pm$ 0.34	88.67 $\pm$ 0.56	3.09

distance under different  $\beta$  values for UGradSL and UGradSL+ are reported in Table 6 and Table 7, respectively.

## E EXPERIMENTS

### E.1 ADDITIONAL EXPERIMENTAL SETTINGS

The datasets and model configurations for the original model training and retraining are given in Table 8. We run all the experiments using PyTorch 1.12 on NVIDIA A5000 GPUs and AMD EPYC 7513 32-Core Processor.

Table 8: The hyperparameters used in the original training and retraining for different models and datasets.

Settings	CIFAR-10			SVHN	CIFAR-100	ImageNet	20 Newsgroups
	ResNet-18	VGG-16	ViT	ResNet-18	ResNet-18	ResNet-18	Bert
Batch Size	256	256	256	256	256	1024	128
Learning rate	$1e^{-2}$	$1e^{-4}$	$1e^{-6}$	$1e^{-2}$	$1e^{-2}$	$1e^{-2}$	$1e^{-4}$
Epochs	160	160	160	160	160	90	60

The settings of the baseline methods are:

- Fine-tuning (FT): FT is to fine-tune the original model  $\theta_o$  trained from  $D_{tr}$  using  $D_r$ . We fix the epoch of FT for 10 epochs for all the datasets except ImageNet. We fine-tune ImageNet for 5 epochs. The learning rate is the same as the original training.
- Fisher forgetting (FF): FF is to perturb the  $\theta_o$  by adding the Gaussian noise, which with a zero mean and a covariance corresponds to the 4th root of the Fisher Information Matrix with respect to (w.r.t.)  $\theta_o$  on  $D_r$  (Golatkar et al., 2020). We perform a greedy search for hyperparameter tuning between  $1e^{-9}$  and  $1e^{-6}$ .
- Influence unlearning (IU): IU uses influence function (Koh & Liang, 2017) to estimate the change from  $\theta_o$  to  $\theta_u$  when one training sample is removed.
- Boundary unlearning<sup>2</sup> (BU): BU unlearns the data by assigning pseudo label and manipulating the decision boundary. It contains boundary shrink and boundary expansion, two types of unlearning methods. The hyper-parameters are the default values in the paper.
- $\ell_1$ -sparse<sup>3</sup>:  $\ell_1$ -sparse improves machine unlearning by integrating the  $\ell_1$  norm-based sparse penalty to the loss function. The learning rate is  $1e^{-3}$  and we search  $\gamma$  in  $[1e^{-5}, 1e^{-1}]$  as given in (Jia et al., 2023).
- SCRUB: SCRUB casts the unlearning problem into a teacher-student framework. We follow the settings exactly the same in the original repo<sup>4</sup> where  $\gamma = 0.99$  and  $\alpha = 0.001$ .
- Random Labeling (RL): Unlike FT, RL is to train the model with the random label rather than the fixed label. The settings are the same as for FT.
- SalUN<sup>5</sup>: SalUN takes the weight saliency into consideration. We search  $\gamma$  from  $[0.5, 0.9]$ .

### E.2 DATASET SPLIT OF DIFFERENT FORGETTING PARADIGMS

We also provide the details of the dataset split for different forgetting paradigms. For classwise forgetting, we remove the whole class from  $D_{tr}$  and  $D_{te}$ . In CIFAR-10 and CIFAR-100, the size of  $D_f$  is 500 and 5000, respectively. For the other datasets, the size of  $D_f$  ranges from the smallest class size to the largest class size because we remove the whole class completely. The selected class to be forgotten is totally random. For random forgetting, we randomly select 10% data from  $D_{tr}$  as  $D_f$ . We make sure the distribution of  $D_f$  is the same as  $D_{tr}$ . For CIFAR-20 in group forgetting, each fine-grained class is in the same size which is 500. The coarse class is 2500.

<sup>2</sup><https://github.com/TY-LEE-KR/Boundary-Unlearning-Code>

<sup>3</sup><https://github.com/OPTML-Group/Unlearn-Sparse>

<sup>4</sup><https://github.com/meghdadk/SCRUB/tree/main>

<sup>5</sup><https://github.com/OPTML-Group/Unlearn-Saliency>

Table 9: The experiment results of class-wise forgetting in 20 Newsgroups and SVHN datasets.

20 Newsgroups	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap (↓)	RTE (↓, min)
Retrain	100.00 $\pm$ 0.00	100.00 $\pm$ 0.00	98.31 $\pm$ 2.56	81.95 $\pm$ 1.69	-	26.25
FT	4.14 $\pm$ 2.11	9.23 $\pm$ 3.40	98.83 $\pm$ 0.86	82.63 $\pm$ 0.73	46.96	1.77
GA	17.12 $\pm$ 9.48	62.03 $\pm$ 5.84	99.99 $\pm$ 0.01	85.41 $\pm$ 0.37	31.50	<b>0.37</b>
IU	0.00 $\pm$ 0.00	0.25 $\pm$ 0.12	100.00 $\pm$ 0.00	85.58 $\pm$ 0.20	51.27	1.52
BS	78.33 $\pm$ 3.47	92.63 $\pm$ 2.19	97.28 $\pm$ 0.99	90.93 $\pm$ 0.81	9.76	1.42
UGradSL	100.00 $\pm$ 0.00	100.00 $\pm$ 0.00	96.31 $\pm$ 4.02	78.54 $\pm$ 5.10	1.35	0.39
UGradSL+	100.00 $\pm$ 0.00	100.00 $\pm$ 0.00	99.76 $\pm$ 0.23	84.21 $\pm$ 0.41	<b>0.93</b>	2.13
SVHN	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap (↓)	RTE (↓, min)
Retrain	100.00 $\pm$ 0.00	100.00 $\pm$ 0.00	100.00 $\pm$ 0.01	95.94 $\pm$ 0.11	-	37.05
FT	6.49 $\pm$ 1.49	99.98 $\pm$ 0.04	100.00 $\pm$ 0.01	96.08 $\pm$ 0.01	23.42	2.42
GA	87.49 $\pm$ 1.94	99.85 $\pm$ 0.09	99.52 $\pm$ 0.03	95.27 $\pm$ 0.21	3.45	<b>0.15</b>
IU	93.55 $\pm$ 2.78	100.00 $\pm$ 0.00	99.54 $\pm$ 0.03	95.64 $\pm$ 0.31	1.80	0.23
BE	85.56 $\pm$ 3.07	99.98 $\pm$ 0.02	99.55 $\pm$ 0.01	95.53 $\pm$ 0.07	3.83	3.17
BS	96.62 $\pm$ 1.14	99.95 $\pm$ 0.09	99.99 $\pm$ 0.00	95.39 $\pm$ 0.18	1.00	3.91
$\ell_1$ -sparse	99.78 $\pm$ 0.31	100.00 $\pm$ 0.00	98.63 $\pm$ 0.01	97.36 $\pm$ 0.18	0.75	2.91
RL	99.99 $\pm$ 0.01	100.00 $\pm$ 0.00	100.00 $\pm$ 0.00	95.44 $\pm$ 0.13	0.13	3.53
EU- $k$	100.00 $\pm$ 0.00	100.00 $\pm$ 0.00	99.61 $\pm$ 0.08	65.56 $\pm$ 2.38	7.69	4.93
CF- $k$	0.09 $\pm$ 0.03	2.18 $\pm$ 2.21	99.34 $\pm$ 0.02	69.87 $\pm$ 4.13	56.12	5.02
SCRUB	99.99 $\pm$ 0.02	100.00 $\pm$ 0.00	100.00 $\pm$ 0.00	95.79 $\pm$ 0.26	<b>0.04</b>	4.97
SalUN	99.74 $\pm$ 0.39	100.00 $\pm$ 0.00	99.53 $\pm$ 0.02	95.00 $\pm$ 1.50	0.42	4.77
UGradSL	90.71 $\pm$ 4.08	99.90 $\pm$ 0.16	99.54 $\pm$ 0.04	95.64 $\pm$ 0.25	2.54	0.23
UGradSL+	100.00 $\pm$ 0.00	100.00 $\pm$ 0.00	99.82 $\pm$ 0.62	94.35 $\pm$ 0.70	0.44	4.56

### E.3 ADDITIONAL CLASS-WISE FORGETTING RESULTS

We present the performance of class-wise forgetting in 20 Newsgroups and SVHN datasets in Table 9. The observation is similar in CIFAR-100 and ImageNet given in Table 1. UGradSL and UGradSL+ can improve the MU performance with acceptable time increment, showing the generalization of the proposed method in different modalities and different dataset sizes.

### E.4 ADDITIONAL RANDOM FORGETTING RESULTS

We present the performance of random forgetting in CIFAR-10 and SVHN datasets in Table 10. The observation is similar in CIFAR-100 and Tiny ImageNet given in Table 2.

### E.5 MU WITH THE OTHER CLASSIFIER

To validate the generalization of the proposed method, we also try the other classification models. We test vision transformer (ViT) and VGG-16 on the task of class-wise forgetting and random forgetting using CIFAR-10, respectively. The results are given in Tables 11 and 12. The observation is similar in Tables 1 and 2, respectively.

Table 11: The experiment results of class-wise forgetting in CIFAR-10 using ViT.

CIFAR-10	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap (↓)	RTE (↓, min)
Retrain	100.00 $\pm$ 0.00	100.00 $\pm$ 0.00	61.41 $\pm$ 0.81	58.94 $\pm$ 1.09	-	189.08
FT	3.97 $\pm$ 0.87	7.60 $\pm$ 1.76	98.29 $\pm$ 0.05	80.44 $\pm$ 0.22	61.70	2.99
GA	33.77 $\pm$ 6.36	40.47 $\pm$ 6.63	89.47 $\pm$ 4.21	71.63 $\pm$ 2.79	41.63	0.32
IU	1.74 $\pm$ 0.09	2.16 $\pm$ 0.61	73.96 $\pm$ 0.01	68.88 $\pm$ 0.00	54.65	0.24
BE	85.56 $\pm$ 3.07	99.98 $\pm$ 0.02	99.55 $\pm$ 0.01	95.53 $\pm$ 0.07	22.30	3.17
UGradSL	68.11 $\pm$ 11.03	73.84 $\pm$ 9.58	84.11 $\pm$ 2.70	68.33 $\pm$ 1.69	22.54	<b>0.22</b>
UGradSL+	99.99 $\pm$ 0.01	99.99 $\pm$ 0.02	94.46 $\pm$ 1.06	77.26 $\pm$ 1.19	<b>12.85</b>	5.86

Table 10: The experiment results of random forgetting in CIFAR-10 and SVHN.

CIFAR-10	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap ( $\downarrow$ )	RTE ( $\downarrow$ , min)
Retrain	8.07 $\pm$ 0.47	17.41 $\pm$ 0.69	100.00 $\pm$ 0.01	91.61 $\pm$ 0.24	-	24.66
FT	1.10 $\pm$ 0.19	4.06 $\pm$ 0.41	99.83 $\pm$ 0.03	93.70 $\pm$ 0.10	5.65	1.58
GA	0.56 $\pm$ 0.01	1.19 $\pm$ 0.05	99.48 $\pm$ 0.02	94.55 $\pm$ 0.05	<b>6.80</b>	<b>0.31</b>
IU	17.51 $\pm$ 2.19	21.39 $\pm$ 1.70	83.28 $\pm$ 2.44	78.13 $\pm$ 2.85	10.91	1.18
BE	0.00 $\pm$ 0.00	0.26 $\pm$ 0.02	100.00 $\pm$ 0.00	95.35 $\pm$ 0.18	7.24	3.17
BS	0.48 $\pm$ 0.07	1.16 $\pm$ 0.04	99.47 $\pm$ 0.01	94.58 $\pm$ 0.03	6.84	1.41
$\ell_1$ -sparse	1.21 $\pm$ 0.38	4.33 $\pm$ 0.52	97.39 $\pm$ 0.31	95.49 $\pm$ 0.18	6.61	1.82
SCRUB	0.70 $\pm$ 0.59	3.88 $\pm$ 1.25	99.59 $\pm$ 0.34	94.22 $\pm$ 0.26	5.98	4.05
Random Label	2.80 $\pm$ 0.37	18.59 $\pm$ 3.48	99.97 $\pm$ 0.01	94.08 $\pm$ 0.12	2.24	1.98
UGradSL	5.87 $\pm$ 0.51	13.33 $\pm$ 0.70	98.82 $\pm$ 0.28	92.17 $\pm$ 0.23	2.01	0.45
UGradSL+	6.03 $\pm$ 0.17	10.65 $\pm$ 0.13	99.79 $\pm$ 0.03	93.64 $\pm$ 0.16	2.76	3.07
UGradSL (Adp)	6.04 $\pm$ 0.11	13.75 $\pm$ 0.32	99.11 $\pm$ 0.01	92.07 $\pm$ 0.02	1.76	1.35
UGradSL+ (Adp)	7.54 $\pm$ 0.43	13.57 $\pm$ 0.12	99.67 $\pm$ 0.00	92.97 $\pm$ 0.17	<b>1.52</b>	9.23
SVHN	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap ( $\downarrow$ )	RTE ( $\downarrow$ , min)
Retrain	4.95 $\pm$ 0.03	15.59 $\pm$ 0.93	99.99 $\pm$ 0.01	95.61 $\pm$ 0.22	-	35.65
FT	0.45 $\pm$ 0.14	2.30 $\pm$ 0.04	99.99 $\pm$ 0.00	95.78 $\pm$ 0.01	4.49	2.76
GA	0.58 $\pm$ 0.04	1.13 $\pm$ 0.02	99.56 $\pm$ 0.01	95.62 $\pm$ 0.01	4.82	0.31
FF	0.45 $\pm$ 0.09	1.30 $\pm$ 0.12	99.55 $\pm$ 0.01	95.49 $\pm$ 0.03	4.84	6.02
BE	0.00 $\pm$ 0.02	0.02 $\pm$ 0.17	100.00 $\pm$ 0.01	96.14 $\pm$ 0.02	5.27	1.03
BS	0.45 $\pm$ 0.14	1.13 $\pm$ 0.05	99.57 $\pm$ 0.03	95.66 $\pm$ 0.01	4.86	4.24
$\ell_1$ -sparse	3.73 $\pm$ 0.78	8.44 $\pm$ 0.34	97.84 $\pm$ 0.28	96.18 $\pm$ 0.33	2.77	<b>0.07</b>
SCRUB	0.35 $\pm$ 0.20	4.96 $\pm$ 0.93	99.94 $\pm$ 0.02	95.36 $\pm$ 0.23	3.88	3.24
RL	8.00 $\pm$ 0.64	29.40 $\pm$ 11.92	98.72 $\pm$ 0.45	94.04 $\pm$ 1.10	4.93	1.79
UGradSL	3.29 $\pm$ 2.53	14.32 $\pm$ 4.56	99.89 $\pm$ 0.02	94.38 $\pm$ 0.28	1.07	0.57
UGradSL+	5.77 $\pm$ 2.93	15.95 $\pm$ 2.26	100.00 $\pm$ 0.00	95.12 $\pm$ 0.50	0.42	4.44
UGradSL (Adp)	3.97 $\pm$ 0.29	14.63 $\pm$ 2.15	99.89 $\pm$ 0.01	94.40 $\pm$ 0.12	0.81	2.20
UGradSL+ (Adp)	5.07 $\pm$ 0.34	15.89 $\pm$ 1.03	100.00 $\pm$ 0.00	95.21 $\pm$ 0.44	<b>0.21</b>	14.33

Table 12: The experiment results of random forgetting across all classes in CIFAR-10 using VGG-16

CIFAR-10	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap ( $\downarrow$ )	RTE ( $\downarrow$ , min)
Retrain	11.41 $\pm$ 0.41	11.97 $\pm$ 0.50	74.65 $\pm$ 0.23	66.13 $\pm$ 0.16	-	9.48
FT	1.32 $\pm$ 0.13	3.48 $\pm$ 0.13	74.24 $\pm$ 0.04	67.04 $\pm$ 0.10	4.98	0.60
GA	1.39 $\pm$ 0.08	2.18 $\pm$ 0.66	73.95 $\pm$ 0.01	66.88 $\pm$ 0.01	5.33	<b>0.14</b>
IU	1.74 $\pm$ 0.09	2.16 $\pm$ 0.61	73.96 $\pm$ 0.01	68.88 $\pm$ 0.00	5.73	0.24
FF	1.39 $\pm$ 0.09	2.21 $\pm$ 0.58	73.95 $\pm$ 0.02	66.87 $\pm$ 0.04	5.32	1.02
BE	0.01 $\pm$ 0.01	0.23 $\pm$ 0.05	99.98 $\pm$ 0.00	94.04 $\pm$ 0.21	19.10	1.09
BS	0.01 $\pm$ 0.01	0.22 $\pm$ 0.03	99.98 $\pm$ 0.01	94.00 $\pm$ 0.14	19.09	3.17
$\ell_1$ -sparse	1.27 $\pm$ 1.13	3.60 $\pm$ 2.41	98.97 $\pm$ 1.13	92.18 $\pm$ 1.46	17.22	0.08
SCRUB	61.16 $\pm$ 50.89	44.65 $\pm$ 43.31	39.26 $\pm$ 50.57	36.95 $\pm$ 46.68	36.75	0.91
UGradSL	13.45 $\pm$ 0.63	11.77 $\pm$ 0.54	65.05 $\pm$ 0.48	58.52 $\pm$ 0.38	<b>4.86</b>	0.19
UGradSL+	12.41 $\pm$ 0.32	14.96 $\pm$ 0.52	65.90 $\pm$ 0.52	58.58 $\pm$ 0.35	5.13	1.08

## E.6 STREISAND EFFECT

From the perspective of security, it is important to make the predicted distributions are almost the same from the forgetting set  $D_f$  and the testing set  $D_{te}$ , which is called Streisand effect. We investigate this effect in the *random forgetting* on CIFAR-10 by plotting confusion matrix as shown in Figure 5. It can be found that our method will not lead to the extra hint of  $D_f$ .

Table 13: Ablation studies of GA ratio  $p$  for random forgetting on CIFAR-10. The forgetting set size is 10% training set. The method is UGradSL. We fix  $\alpha$  as -0.4. The first row is the retraining results for reference.

$p$	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap ( $\downarrow$ )
-	8.07	17.41	100.00	91.61	-
0.80	12.47 $\pm$ 1.01	20.24 $\pm$ 2.12	94.11 $\pm$ 0.71	88.21 $\pm$ 0.57	4.13
0.81	11.57 $\pm$ 1.69	19.68 $\pm$ 3.31	94.39 $\pm$ 1.41	88.56 $\pm$ 1.04	3.61
0.82	10.61 $\pm$ 0.08	17.85 $\pm$ 0.97	94.92 $\pm$ 0.18	88.94 $\pm$ 0.37	2.68
0.83	9.64 $\pm$ 0.52	16.49 $\pm$ 0.96	95.54 $\pm$ 0.71	89.63 $\pm$ 0.71	2.23
0.84	9.33 $\pm$ 1.04	15.84 $\pm$ 1.17	95.43 $\pm$ 0.95	89.48 $\pm$ 0.71	2.38
0.85	8.27 $\pm$ 0.82	14.74 $\pm$ 1.36	96.05 $\pm$ 0.41	90.08 $\pm$ 0.40	2.09
0.86	7.96 $\pm$ 0.42	15.45 $\pm$ 2.00	96.10 $\pm$ 0.44	90.22 $\pm$ 0.18	<b>1.84</b>
0.87	7.51 $\pm$ 0.26	15.26 $\pm$ 3.47	96.05 $\pm$ 0.18	90.20 $\pm$ 0.51	2.02
0.88	6.87 $\pm$ 0.37	13.18 $\pm$ 1.26	96.43 $\pm$ 0.35	90.29 $\pm$ 0.64	2.58
0.89	6.91 $\pm$ 0.56	14.44 $\pm$ 2.46	96.38 $\pm$ 0.73	90.47 $\pm$ 0.36	2.22
0.90	6.92 $\pm$ 1.08	13.60 $\pm$ 3.42	96.00 $\pm$ 0.50	90.26 $\pm$ 0.14	2.58
0.91	6.44 $\pm$ 1.30	14.16 $\pm$ 2.27	95.93 $\pm$ 1.18	90.17 $\pm$ 0.72	2.60
0.92	6.50 $\pm$ 0.69	14.35 $\pm$ 0.72	95.64 $\pm$ 0.50	90.06 $\pm$ 0.12	2.64
0.93	5.88 $\pm$ 0.82	14.84 $\pm$ 1.26	96.03 $\pm$ 0.84	90.31 $\pm$ 0.54	2.51
0.94	5.65 $\pm$ 0.30	13.55 $\pm$ 0.78	96.25 $\pm$ 0.44	90.54 $\pm$ 0.10	2.77
0.95	6.13 $\pm$ 1.29	13.14 $\pm$ 2.43	95.73 $\pm$ 1.03	89.88 $\pm$ 0.75	3.05
0.96	6.07 $\pm$ 0.91	14.28 $\pm$ 2.15	95.64 $\pm$ 0.79	90.15 $\pm$ 0.36	2.74
0.97	5.83 $\pm$ 1.25	14.07 $\pm$ 1.98	95.20 $\pm$ 0.98	89.67 $\pm$ 0.59	3.08
0.98	5.73 $\pm$ 0.84	13.19 $\pm$ 1.99	95.43 $\pm$ 0.98	89.82 $\pm$ 0.38	3.23
0.99	5.83 $\pm$ 1.05	12.98 $\pm$ 1.37	94.99 $\pm$ 0.79	89.46 $\pm$ 0.59	3.46

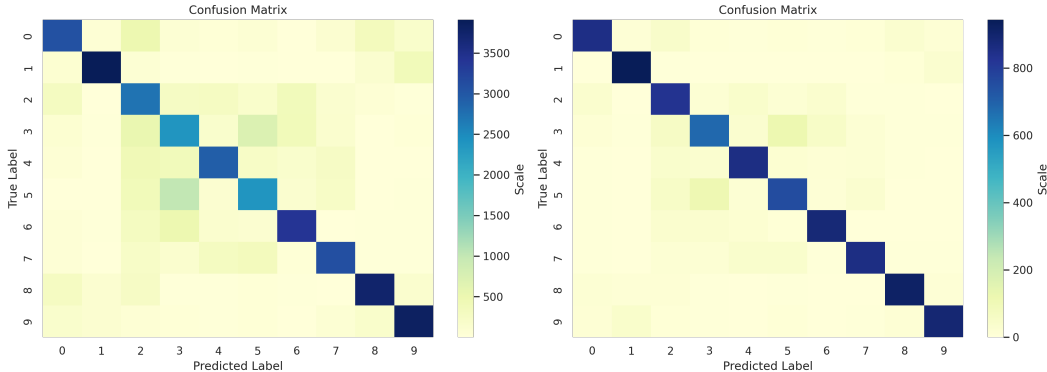


Figure 5: The confusion matrix of testing set and forgetting set  $D_f$  using our method on CIFAR-10 with random forgetting across all the classes. There is no big difference between the prediction distribution. Our method will not make  $D_f$  more distinguishable.

## E.7 ABLATION STUDY: FORGETTING SET SIZE

Since the size of the forgetting set can affect unlearning performance, we further evaluate the robustness of our method under varying forgetting ratios. In addition to the 10% random forgetting results reported in Table 2 and Table 4, we consider forgetting set sizes of 20%, 30%, 40%, and 50% of the training data on CIFAR-10 and CIFAR-100. The results are summarized in Tables 17 and 18.

## E.8 ABLATION STUDY: GA RATIO $p$

In addition to an overview of the performance fluctuation in Figure 3. We provide the specific value of the ablation study regarding GA ratio  $p$ . We test the performance on random forgetting on CIFAR-10. The forgetting set size is 10% of the training set. The results of UGradSL and UGradSL+ are given in Tables 13 and 14.

Table 14: Ablation studies of GA ratio  $p$  for random forgetting on CIFAR-10. The forgetting set size is 10% training set. The method is UGradSL+. We fix  $\alpha$  as -0.4. The first row is the retraining results for reference.

$p$	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap (↓)
-	8.07	17.41	100.00	91.61	-
0.80	18.84 $\pm$ 0.71	26.78 $\pm$ 1.78	91.95 $\pm$ 0.86	86.05 $\pm$ 1.21	8.44
0.81	17.00 $\pm$ 0.47	24.55 $\pm$ 0.76	93.21 $\pm$ 0.21	87.50 $\pm$ 0.49	6.74
0.82	16.45 $\pm$ 0.33	24.22 $\pm$ 0.46	93.60 $\pm$ 0.49	87.64 $\pm$ 0.25	6.39
0.83	14.45 $\pm$ 0.73	21.73 $\pm$ 0.82	94.66 $\pm$ 0.38	88.59 $\pm$ 0.19	4.76
0.84	13.44 $\pm$ 0.77	20.92 $\pm$ 1.15	94.67 $\pm$ 0.71	88.67 $\pm$ 0.44	4.29
0.85	12.57 $\pm$ 0.65	19.18 $\pm$ 0.64	95.29 $\pm$ 1.02	89.33 $\pm$ 1.09	3.32
0.86	11.42 $\pm$ 0.14	18.34 $\pm$ 0.61	95.56 $\pm$ 0.46	89.49 $\pm$ 0.27	2.71
0.87	10.90 $\pm$ 0.72	17.22 $\pm$ 0.51	95.79 $\pm$ 0.39	89.77 $\pm$ 0.81	2.27
0.88	10.13 $\pm$ 0.42	17.85 $\pm$ 2.11	95.97 $\pm$ 0.17	90.03 $\pm$ 0.60	2.03
0.89	8.98 $\pm$ 0.29	14.94 $\pm$ 0.09	96.20 $\pm$ 0.27	90.23 $\pm$ 0.33	2.14
0.90	8.41 $\pm$ 0.33	16.87 $\pm$ 1.17	96.53 $\pm$ 0.03	90.64 $\pm$ 0.09	1.33
0.91	8.01 $\pm$ 0.30	17.33 $\pm$ 1.17	96.50 $\pm$ 0.36	90.68 $\pm$ 0.40	<b>1.14</b>
0.92	7.74 $\pm$ 0.33	15.62 $\pm$ 1.80	96.28 $\pm$ 0.26	90.48 $\pm$ 0.46	1.75
0.93	6.67 $\pm$ 0.12	15.93 $\pm$ 0.22	96.80 $\pm$ 0.10	90.96 $\pm$ 0.34	1.67
0.94	6.79 $\pm$ 0.71	16.47 $\pm$ 0.52	96.42 $\pm$ 0.83	90.74 $\pm$ 0.45	1.67
0.95	6.03 $\pm$ 0.26	14.82 $\pm$ 1.39	96.76 $\pm$ 0.41	90.94 $\pm$ 0.35	2.14
0.96	5.78 $\pm$ 0.24	14.79 $\pm$ 1.14	96.90 $\pm$ 0.19	91.30 $\pm$ 0.16	2.08
0.97	5.98 $\pm$ 0.49	14.96 $\pm$ 0.34	96.56 $\pm$ 0.45	90.81 $\pm$ 0.53	2.20
0.98	6.46 $\pm$ 0.74	15.15 $\pm$ 1.76	95.52 $\pm$ 0.67	90.15 $\pm$ 0.89	2.45
0.99	5.67 $\pm$ 0.27	14.40 $\pm$ 1.18	96.17 $\pm$ 0.46	90.61 $\pm$ 0.24	2.56

Table 15: The ablation study of smoothing rate  $\alpha$  for random forgetting on CIFAR-10. The forgetting set size is 10% training set. The method we use is UGradSL. We fix  $p$  as 0.9.

$\alpha$	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap (↓)
-	8.07	17.41	100.00	91.61	-
-0.9	8.17 $\pm$ 1.74	14.96 $\pm$ 2.47	95.81 $\pm$ 1.97	89.97 $\pm$ 1.46	<b>2.10</b>
-0.8	6.98 $\pm$ 0.47	13.41 $\pm$ 0.99	96.67 $\pm$ 0.88	90.75 $\pm$ 0.22	2.32
-0.7	7.23 $\pm$ 0.56	14.33 $\pm$ 1.23	96.28 $\pm$ 0.11	90.47 $\pm$ 0.26	2.20
-0.6	6.69 $\pm$ 0.22	12.93 $\pm$ 0.67	96.46 $\pm$ 0.39	90.64 $\pm$ 0.04	2.59
-0.5	6.56 $\pm$ 0.29	13.00 $\pm$ 0.50	96.58 $\pm$ 0.23	90.66 $\pm$ 0.20	2.57
-0.4	6.92 $\pm$ 1.08	13.60 $\pm$ 3.42	96.00 $\pm$ 0.50	90.26 $\pm$ 0.14	2.58
-0.3	6.32 $\pm$ 0.43	13.63 $\pm$ 0.67	96.18 $\pm$ 0.41	90.52 $\pm$ 0.27	2.61
-0.2	6.95 $\pm$ 0.54	13.98 $\pm$ 1.99	95.41 $\pm$ 0.68	89.65 $\pm$ 0.56	2.77
-0.1	7.13 $\pm$ 1.44	14.47 $\pm$ 1.91	95.08 $\pm$ 1.55	89.57 $\pm$ 1.04	2.71

## E.9 SMOOTHING RATIO $\alpha$

Similar to  $p$ , we report the detailed results regarding the smoothing rate  $\alpha$ . The results of UGradSL and UGradSL+ are given in Tables 15 and 16.

## E.10 GRADIENT ANALYSIS

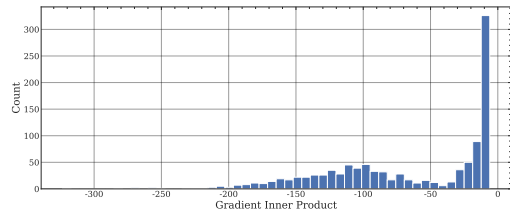
As mentioned in Section 3.3,  $\langle \Delta\theta_r - \Delta\theta_f, \Delta\theta_n - \Delta\theta_f \rangle \leq 0$  is always practically valid. We check the results on CelebA dataset (ResNet-18), ImageNet (ViT), CIFAR-100 (VGG-16) and CIFAR-10 (ResNet-18). The distribution of  $\langle \Delta\theta_r - \Delta\theta_f, \Delta\theta_n - \Delta\theta_f \rangle$  is shown in Figure 6, which aligns with our assumption.

## E.11 THE DIFFERENCE BETWEEN UGRADSL AND UGRADSL+

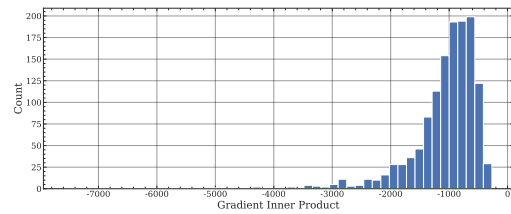
Although UGradSL and UGradSL+ look similar, the intuition of these two method is totally different because of the difference between FT and GA. We conducted experiments to illustrate the difference between GA and FT as well as UGradSL and UGradSL+. The results are given in Table 19. The dataset and forgetting paradigm is CIFAR-10 random forgetting. It can be found that the difference becomes much larger when the number of epochs is over 8. When the number of epochs is 10, the model becomes unusable because TA is less than 10%. We also report the performance of UGradSL

Table 16: The ablation study of smoothing rate  $\alpha$  for random forgetting on CIFAR-10. The forgetting set size is 10% training set. The method we use is UGradSL+. We fix  $p$  as 0.9.

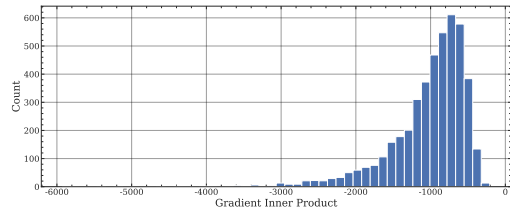
$\alpha$	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap ( $\downarrow$ )
-	8.07	17.41	100.00	91.61	-
-0.9	11.59 $\pm$ 0.40	19.41 $\pm$ 0.59	95.77 $\pm$ 0.58	89.47 $\pm$ 0.56	2.97
-0.8	10.68 $\pm$ 0.27	18.41 $\pm$ 0.48	95.94 $\pm$ 0.24	89.91 $\pm$ 0.36	2.35
-0.7	10.12 $\pm$ 1.01	16.88 $\pm$ 0.69	96.07 $\pm$ 0.89	90.08 $\pm$ 0.78	2.01
-0.6	8.98 $\pm$ 0.15	16.29 $\pm$ 0.87	96.64 $\pm$ 0.19	90.75 $\pm$ 0.05	1.56
-0.5	9.07 $\pm$ 0.21	15.83 $\pm$ 0.25	96.65 $\pm$ 0.34	90.43 $\pm$ 0.50	1.78
-0.4	8.41 $\pm$ 0.33	16.87 $\pm$ 1.17	96.53 $\pm$ 0.03	90.64 $\pm$ 0.09	<b>1.33</b>
-0.3	8.59 $\pm$ 0.07	16.86 $\pm$ 2.15	96.24 $\pm$ 0.38	90.20 $\pm$ 0.15	1.56
-0.2	7.55 $\pm$ 0.18	16.68 $\pm$ 1.60	96.43 $\pm$ 0.14	90.86 $\pm$ 0.33	1.39
-0.1	7.57 $\pm$ 0.18	17.32 $\pm$ 0.23	96.15 $\pm$ 0.38	90.34 $\pm$ 0.27	1.43



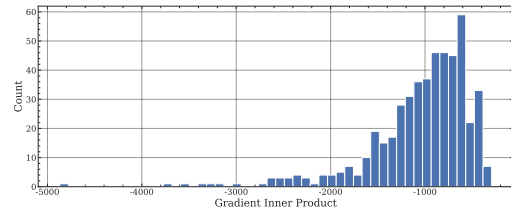
(a) The group forgetting on CelebA using ResNet-18



(b) The class-wise forgetting on ImageNet using ViT



(c) The random forgetting on CIFAR-100 using VGG-16



(d) The class-wise forgetting on CIFAR-10 using ResNet-18

Figure 6: The distribution of  $\langle \Delta\theta_r - \Delta\theta_f, \Delta\theta_n - \Delta\theta_f \rangle$  using multiple models on multiple datasets.

and UGradSL+ in different epochs. For UGradSL, when the epochs are over 14, the model cannot be used at all. For UGradSL+, the algorithm is much more stable, showing the very good adaptive capability.

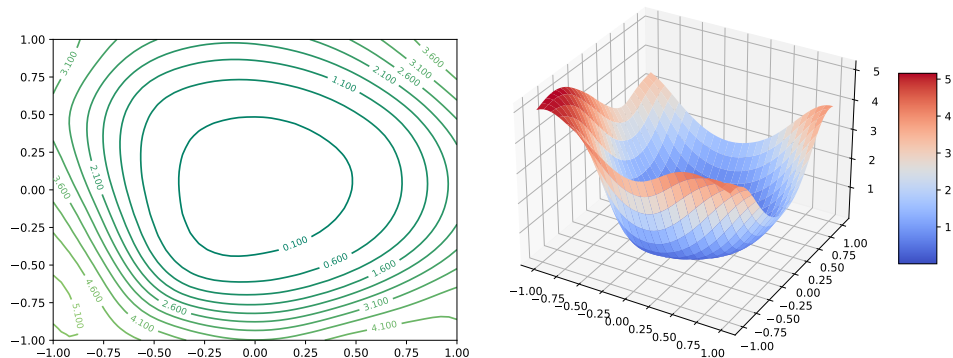


Figure 7: The loss landscape of  $\theta_r$  on CIFAR-10 and the model is ResNet-18.

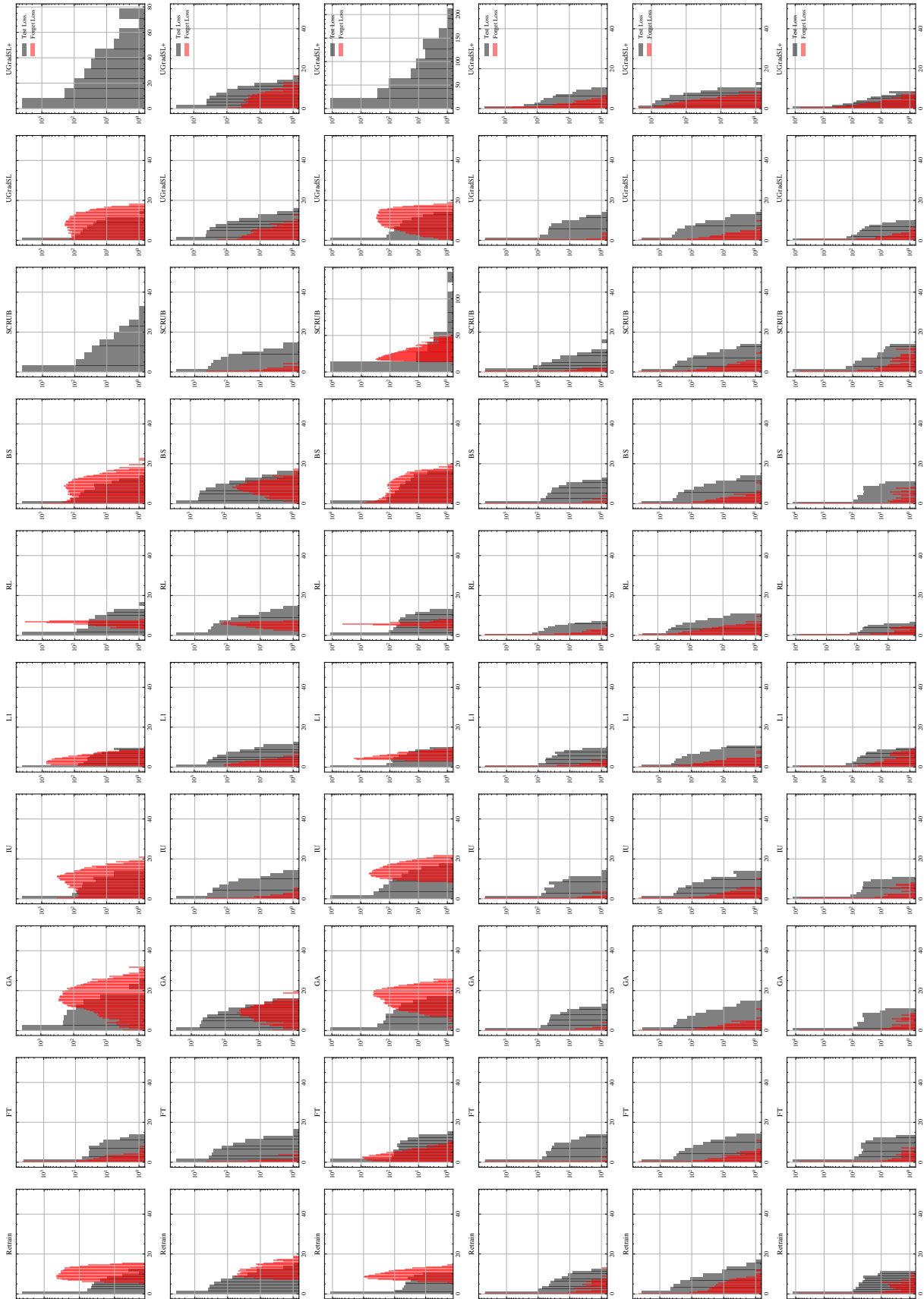


Figure 8: The distributions of the cross-entropy losses for the forget and test instances from the unlearned models. The y-axis is in log scale for better visualization. From the first to the last figure, they are random forgetting on CIFAR-10, SVHN and class-wise forgetting on CIFAR-10, CIFAR-100, SVHN.

Table 17: MU Performance across different forgetting data amounts on ResNet-18, pre-trained on CIFAR-10 dataset, for random data forgetting.

Method	Random Set Size (10%)					
	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap (↓)	RTE (↓, min)
Retrain	8.07	17.41	100.00	91.61	-	24.66
FT	1.10 $\pm$ 0.19	4.06 $\pm$ 0.40	98.83 $\pm$ 0.03	93.70 $\pm$ 0.10	5.90	1.58
RL	6.39 $\pm$ 1.09	0.00 $\pm$ 0.00	99.50 $\pm$ 0.10	99.04 $\pm$ 0.08	6.76	1.92
GA	0.56 $\pm$ 0.01	1.19 $\pm$ 0.05	99.48 $\pm$ 0.02	94.55 $\pm$ 0.05	6.80	<b>0.31</b>
IU	17.51 $\pm$ 2.19	21.39 $\pm$ 1.70	98.00 $\pm$ 0.38	98.11 $\pm$ 0.38	5.48	1.18
BE	0.00 $\pm$ 0.00	0.26 $\pm$ 0.02	100.00 $\pm$ 0.00	95.35 $\pm$ 0.18	7.24	1.37
BS	0.37 $\pm$ 0.10	1.10 $\pm$ 0.43	99.93 $\pm$ 0.01	98.97 $\pm$ 0.02	7.86	1.21
$\ell_1$ -sparse	2.80 $\pm$ 0.37	19.59 $\pm$ 3.48	99.07 $\pm$ 0.04	98.00 $\pm$ 0.12	3.69	1.98
SalUn	46.95 $\pm$ 0.15	86.33 $\pm$ 2.58	97.75 $\pm$ 0.42	97.22 $\pm$ 0.77	28.92	2.42
UGradSL	5.87 $\pm$ 0.50	13.33 $\pm$ 0.20	98.82 $\pm$ 0.28	92.17 $\pm$ 0.20	<b>2.01</b>	0.45
UGradSL+	6.03 $\pm$ 0.17	10.65 $\pm$ 0.13	99.79 $\pm$ 0.03	93.64 $\pm$ 0.16	2.76	3.07
Method	Random Set Size (20%)					
	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap (↓)	RTE (↓, min)
Retrain	5.31	13.30	100.00	94.10	-	38.74
FT	0.76 $\pm$ 4.55	2.69 $\pm$ 10.61	99.89 $\pm$ 0.11	93.97 $\pm$ 0.13	3.85	2.17
RL	6.47 $\pm$ 1.16	28.62 $\pm$ 15.32	99.60 $\pm$ 0.40	92.39 $\pm$ 1.71	4.65	2.65
GA	0.67 $\pm$ 4.64	1.44 $\pm$ 11.86	99.48 $\pm$ 0.52	94.42 $\pm$ 0.32	4.33	0.26
IU	2.91 $\pm$ 2.40	5.53 $\pm$ 7.77	97.30 $\pm$ 2.70	90.64 $\pm$ 3.46	4.08	3.29
BE	0.57 $\pm$ 4.74	1.64 $\pm$ 11.66	99.44 $\pm$ 0.56	94.32 $\pm$ 0.22	4.29	0.53
BS	0.62 $\pm$ 4.69	1.62 $\pm$ 11.68	99.46 $\pm$ 0.54	94.20 $\pm$ 0.10	4.25	0.86
$\ell_1$ -sparse	3.92 $\pm$ 1.39	8.94 $\pm$ 4.36	98.09 $\pm$ 1.91	91.92 $\pm$ 2.18	2.46	2.20
SalUn	3.73 $\pm$ 1.58	13.18 $\pm$ 0.12	98.61 $\pm$ 1.39	92.75 $\pm$ 1.35	1.11	2.66
UGradSL	6.07 $\pm$ 0.70	13.82 $\pm$ 1.03	95.71 $\pm$ 0.17	90.19 $\pm$ 0.23	2.37	<b>0.24</b>
UGradSL+	6.39 $\pm$ 0.19	12.34 $\pm$ 1.79	97.08 $\pm$ 0.44	90.91 $\pm$ 0.95	<b>2.04</b>	0.31
Method	Random Set Size (30%)					
	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap (↓)	RTE (↓, min)
Retrain	6.64	14.60	100.00	92.78	-	33.65
FT	0.56 $\pm$ 6.08	1.66 $\pm$ 12.94	99.83 $\pm$ 0.17	94.22 $\pm$ 1.44	5.16	1.98
RL	6.89 $\pm$ 0.25	31.09 $\pm$ 16.49	99.36 $\pm$ 0.64	91.35 $\pm$ 1.43	4.70	2.63
GA	0.65 $\pm$ 5.99	1.50 $\pm$ 13.10	99.46 $\pm$ 0.54	94.44 $\pm$ 1.66	5.32	2.40
IU	3.95 $\pm$ 2.69	7.26 $\pm$ 7.34	96.22 $\pm$ 3.78	89.61 $\pm$ 3.17	4.24	3.32
BE	0.63 $\pm$ 6.01	3.35 $\pm$ 11.25	99.39 $\pm$ 0.61	94.19 $\pm$ 1.41	4.82	0.81
BS	0.63 $\pm$ 6.01	2.88 $\pm$ 11.72	99.39 $\pm$ 0.61	94.15 $\pm$ 1.37	4.93	1.28
$\ell_1$ -sparse	4.70 $\pm$ 1.94	9.97 $\pm$ 4.63	97.63 $\pm$ 2.37	91.19 $\pm$ 1.59	2.63	1.99
SalUn	6.22 $\pm$ 0.42	14.11 $\pm$ 0.49	95.91 $\pm$ 4.09	90.72 $\pm$ 2.06	1.76	2.64
UGradSL	6.78 $\pm$ 0.66	15.96 $\pm$ 0.12	96.94 $\pm$ 0.56	90.72 $\pm$ 0.80	1.66	<b>0.70</b>
UGradSL+	6.36 $\pm$ 0.65	14.99 $\pm$ 0.82	97.35 $\pm$ 0.79	91.10 $\pm$ 1.10	<b>1.25</b>	0.53
Method	Random Set Size (40%)					
	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap (↓)	RTE (↓, min)
Retrain	7.01	18.37	100.00	92.52	-	28.47
FT	0.77 $\pm$ 6.24	2.88 $\pm$ 15.49	99.96 $\pm$ 0.04	94.27 $\pm$ 1.75	5.88	1.62
RL	5.02 $\pm$ 1.99	37.76 $\pm$ 19.39	99.61 $\pm$ 0.39	92.14 $\pm$ 0.38	5.54	2.68
GA	0.67 $\pm$ 6.34	1.57 $\pm$ 16.80	99.47 $\pm$ 0.53	94.38 $\pm$ 1.86	6.38	<b>0.53</b>
IU	7.89 $\pm$ 0.88	10.99 $\pm$ 7.38	92.21 $\pm$ 7.79	86.15 $\pm$ 6.37	5.60	3.27
BE	0.86 $\pm$ 6.15	15.72 $\pm$ 2.65	99.27 $\pm$ 0.73	93.46 $\pm$ 0.94	2.62	1.04
BS	1.18 $\pm$ 5.83	13.97 $\pm$ 4.40	98.94 $\pm$ 1.06	93.01 $\pm$ 0.49	2.95	1.72
$\ell_1$ -sparse	2.84 $\pm$ 4.17	7.09 $\pm$ 11.28	98.75 $\pm$ 1.25	92.20 $\pm$ 0.32	4.26	1.63
SalUn	6.86 $\pm$ 0.15	15.15 $\pm$ 3.22	95.01 $\pm$ 4.99	89.76 $\pm$ 2.76	2.78	2.67
UGradSL	5.81 $\pm$ 0.11	14.98 $\pm$ 2.65	97.31 $\pm$ 1.06	90.73 $\pm$ 0.48	<b>2.27</b>	0.62
UGradSL+	5.82 $\pm$ 0.37	14.53 $\pm$ 1.83	97.11 $\pm$ 0.40	90.74 $\pm$ 0.38	2.42	0.63
Method	Random Set Size (50%)					
	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap (↓)	RTE (↓, min)
Retrain	7.91	19.29	100.00	91.72	-	23.90
FT	0.44 $\pm$ 7.47	2.15 $\pm$ 17.14	99.96 $\pm$ 0.04	94.23 $\pm$ 2.51	6.79	1.31
RL	7.61 $\pm$ 0.30	37.36 $\pm$ 18.07	99.67 $\pm$ 0.33	92.83 $\pm$ 1.11	4.95	2.65
GA	0.40 $\pm$ 7.51	1.22 $\pm$ 18.07	99.61 $\pm$ 0.39	94.34 $\pm$ 2.62	7.15	<b>0.66</b>
IU	3.97 $\pm$ 3.94	7.29 $\pm$ 12.00	96.21 $\pm$ 3.79	90.00 $\pm$ 1.72	5.36	3.25
BE	3.08 $\pm$ 4.83	24.87 $\pm$ 5.58	96.84 $\pm$ 3.16	90.41 $\pm$ 1.31	3.72	1.31
BS	9.76 $\pm$ 1.85	32.15 $\pm$ 12.86	90.19 $\pm$ 9.81	83.71 $\pm$ 8.01	8.13	2.12
$\ell_1$ -sparse	1.44 $\pm$ 6.47	4.76 $\pm$ 14.53	99.52 $\pm$ 0.48	93.13 $\pm$ 1.41	5.72	1.31
SalUn	7.75 $\pm$ 0.16	16.99 $\pm$ 2.30	94.28 $\pm$ 5.72	89.29 $\pm$ 2.43	2.65	2.68
UGradSL	6.83 $\pm$ 0.23	12.73 $\pm$ 1.66	97.62 $\pm$ 0.71	90.27 $\pm$ 0.55	2.87	0.77
UGradSL+	6.13 $\pm$ 1.35	16.49 $\pm$ 2.73	97.84 $\pm$ 0.34	90.84 $\pm$ 0.69	<b>1.91</b>	0.77

Table 18: MU Performance across different forgetting data amounts on ResNet-18, pre-trained on CIFAR-100 dataset, for random data forgetting.

Method	Random Set Size (10%)					
	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap (↓)	RTE (↓, min)
Retrain	29.47	53.50	99.98	70.51	-	25.01
FT	2.55 $\pm$ 0.03	10.59 $\pm$ 0.27	99.95 $\pm$ 0.01	75.95 $\pm$ 0.05	18.83	1.95
RL	4.06 $\pm$ 0.37	50.12 $\pm$ 3.48	99.92 $\pm$ 0.02	71.30 $\pm$ 0.36	7.41	1.20
GA	2.58 $\pm$ 0.06	5.95 $\pm$ 0.17	97.45 $\pm$ 0.02	76.09 $\pm$ 0.01	20.64	0.29
IU	15.71 $\pm$ 5.19	18.69 $\pm$ 4.12	84.65 $\pm$ 5.19	62.20 $\pm$ 4.17	18.05	1.20
BE	0.01 $\pm$ 0.00	1.45 $\pm$ 0.02	98.22 $\pm$ 1.26	78.26 $\pm$ 0.00	22.76	<b>0.24</b>
BS	2.20 $\pm$ 2.11	10.73 $\pm$ 9.37	98.22 $\pm$ 1.26	70.23 $\pm$ 1.67	18.02	0.34
$\ell_1$ -sparse	8.19 $\pm$ 0.38	19.11 $\pm$ 0.52	88.39 $\pm$ 0.31	80.26 $\pm$ 0.16	19.25	1.00
SalUn	35.23 $\pm$ 0.32	89.39 $\pm$ 0.46	99.53 $\pm$ 0.04	64.26 $\pm$ 0.58	12.10	3.33
UGradSL	18.36 $\pm$ 0.17	40.71 $\pm$ 0.13	98.38 $\pm$ 0.03	68.23 $\pm$ 0.16	6.95	0.55
UGradSL+	21.69 $\pm$ 0.59	49.47 $\pm$ 1.25	99.87 $\pm$ 0.34	73.60 $\pm$ 0.26	<b>3.75</b>	3.52
Method	Random Set Size (20%)					
	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap (↓)	RTE (↓, min)
Retrain	26.84	52.41	99.99	73.88	-	36.88
FT	2.70 $\pm$ 24.14	11.63 $\pm$ 40.78	99.95 $\pm$ 0.04	75.51 $\pm$ 1.63	16.65	2.05
RL	54.74 $\pm$ 27.90	97.32 $\pm$ 44.91	99.47 $\pm$ 0.52	65.59 $\pm$ 8.29	20.41	2.11
GA	6.79 $\pm$ 20.05	13.22 $\pm$ 39.19	94.11 $\pm$ 5.88	71.39 $\pm$ 2.49	16.90	<b>0.26</b>
IU	5.34 $\pm$ 21.50	11.79 $\pm$ 40.62	95.54 $\pm$ 4.45	70.89 $\pm$ 2.99	17.39	3.77
BE	2.51 $\pm$ 24.33	6.70 $\pm$ 45.71	97.38 $\pm$ 2.61	75.07 $\pm$ 1.19	18.46	0.49
BS	2.53 $\pm$ 24.31	6.57 $\pm$ 45.84	97.38 $\pm$ 2.61	75.05 $\pm$ 1.17	18.48	0.82
$\ell_1$ -sparse	37.83 $\pm$ 10.99	38.90 $\pm$ 13.51	76.63 $\pm$ 23.36	58.79 $\pm$ 15.09	15.74	2.05
SalUn	25.83 $\pm$ 1.01	64.69 $\pm$ 12.28	96.01 $\pm$ 3.98	65.87 $\pm$ 8.01	6.32	2.12
UGradSL	30.10 $\pm$ 1.03	47.39 $\pm$ 1.17	93.49 $\pm$ 0.24	64.99 $\pm$ 0.04	<b>5.92</b>	0.83
UGradSL+	27.29 $\pm$ 0.99	35.92 $\pm$ 0.94	93.36 $\pm$ 0.03	66.59 $\pm$ 0.37	7.71	0.59
Method	Random Set Size (30%)					
	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap (↓)	RTE (↓, min)
Retrain	28.52	52.24	99.98	70.91	-	32.92
FT	2.65 $\pm$ 25.87	11.18 $\pm$ 41.06	99.94 $\pm$ 0.04	75.17 $\pm$ 4.26	17.81	1.44
RL	51.46 $\pm$ 22.94	96.34 $\pm$ 44.10	99.32 $\pm$ 0.66	62.77 $\pm$ 8.14	18.96	2.14
GA	2.40 $\pm$ 26.12	5.70 $\pm$ 46.54	97.39 $\pm$ 2.59	75.33 $\pm$ 4.42	19.92	<b>0.40</b>
IU	5.96 $\pm$ 22.56	12.63 $\pm$ 39.61	94.59 $\pm$ 5.39	69.74 $\pm$ 1.17	17.18	3.76
BE	2.44 $\pm$ 26.08	6.53 $\pm$ 45.71	97.37 $\pm$ 2.61	74.77 $\pm$ 3.86	19.56	0.76
BS	2.49 $\pm$ 26.03	6.40 $\pm$ 45.84	97.33 $\pm$ 2.65	74.65 $\pm$ 3.74	19.56	1.24
$\ell_1$ -sparse	38.45 $\pm$ 9.93	38.52 $\pm$ 13.72	76.36 $\pm$ 23.62	58.09 $\pm$ 12.82	15.02	1.47
SalUn	27.34 $\pm$ 1.18	62.99 $\pm$ 10.75	94.50 $\pm$ 5.48	63.10 $\pm$ 7.81	6.31	2.16
UGradSL	30.10 $\pm$ 0.12	47.39 $\pm$ 2.08	93.49 $\pm$ 0.74	64.99 $\pm$ 1.53	<b>4.71</b>	0.83
UGradSL+	24.89 $\pm$ 0.24	44.60 $\pm$ 0.94	94.90 $\pm$ 0.88	66.16 $\pm$ 0.78	5.28	0.79
Method	Random Set Size (40%)					
	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap (↓)	RTE (↓, min)
Retrain	30.07	58.06	99.99	69.87	-	28.29
FT	2.66 $\pm$ 27.41	11.05 $\pm$ 47.01	99.95 $\pm$ 0.04	75.35 $\pm$ 5.48	19.99	1.51
RL	51.75 $\pm$ 21.68	95.78 $\pm$ 37.72	99.27 $\pm$ 0.72	59.41 $\pm$ 10.46	17.64	2.12
GA	2.46 $\pm$ 27.61	5.91 $\pm$ 52.15	97.39 $\pm$ 2.60	75.40 $\pm$ 5.53	21.97	<b>0.51</b>
IU	4.58 $\pm$ 25.49	10.32 $\pm$ 47.74	96.29 $\pm$ 3.70	70.92 $\pm$ 1.05	19.49	3.78
BE	2.54 $\pm$ 27.53	7.44 $\pm$ 50.62	97.35 $\pm$ 2.64	74.56 $\pm$ 4.69	21.37	1.00
BS	2.70 $\pm$ 27.37	7.63 $\pm$ 50.43	97.26 $\pm$ 2.73	74.10 $\pm$ 4.23	21.19	1.66
$\ell_1$ -sparse	38.49 $\pm$ 8.42	40.21 $\pm$ 17.85	78.43 $\pm$ 21.56	57.66 $\pm$ 12.21	15.01	1.52
SalUn	25.54 $\pm$ 4.53	60.08 $\pm$ 2.02	94.64 $\pm$ 5.35	62.52 $\pm$ 7.35	4.81	2.14
UGradSL	30.07 $\pm$ 1.58	49.23 $\pm$ 1.07	95.30 $\pm$ 0.34	64.52 $\pm$ 0.28	<b>4.72</b>	1.08
UGradSL+	30.42 $\pm$ 0.77	45.94 $\pm$ 1.41	93.98 $\pm$ 0.50	63.21 $\pm$ 0.35	6.29	0.77
Method	Random Set Size (50%)					
	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap (↓)	RTE (↓, min)
Retrain	32.69	61.15	99.99	67.22	-	25.01
FT	2.71 $\pm$ 29.98	10.71 $\pm$ 50.44	99.96 $\pm$ 0.03	75.11 $\pm$ 7.89	22.08	1.25
RL	50.52 $\pm$ 17.83	95.91 $\pm$ 34.76	99.47 $\pm$ 0.52	56.75 $\pm$ 10.47	15.90	2.13
GA	2.61 $\pm$ 30.08	5.92 $\pm$ 55.23	97.49 $\pm$ 2.50	75.27 $\pm$ 8.05	23.97	<b>0.66</b>
IU	12.64 $\pm$ 20.05	17.54 $\pm$ 43.61	87.96 $\pm$ 12.03	62.76 $\pm$ 4.46	20.04	3.80
BE	2.76 $\pm$ 29.93	8.85 $\pm$ 52.30	97.39 $\pm$ 2.60	74.05 $\pm$ 6.83	22.92	1.26
BS	2.99 $\pm$ 29.70	8.76 $\pm$ 52.39	97.24 $\pm$ 2.75	73.38 $\pm$ 6.16	22.75	2.08
$\ell_1$ -sparse	39.86 $\pm$ 7.17	40.43 $\pm$ 20.72	78.17 $\pm$ 21.82	55.65 $\pm$ 11.57	15.32	1.26
SalUn	26.17 $\pm$ 6.52	59.47 $\pm$ 1.68	94.04 $\pm$ 5.95	61.39 $\pm$ 5.83	5.00	2.13
UGradSL	33.80 $\pm$ 1.61	57.38 $\pm$ 2.31	95.29 $\pm$ 0.11	56.88 $\pm$ 0.80	<b>4.98</b>	0.95
UGradSL+	32.20 $\pm$ 0.49	49.20 $\pm$ 1.44	94.47 $\pm$ 0.69	61.53 $\pm$ 0.97	5.91	0.75

Table 19: The difference between GA and FT as well as UGradSL and UGradSL+ on CIFAR-10 regarding the number of epochs. The forgetting paradigm is random forgetting.

Epoch	Gradient Ascent					Fine-tuning				
	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap (↓)	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap (↓)
5	0	0.32	95.31	100	9.56	0.04	0.34	95.13	99.99	9.59
6	0	0.40	95.34	100	9.53	-	-	-	-	-
7	0.82	2.22	93.24	99.26	9.21	-	-	-	-	-
8	3.44	4.78	90.80	96.18	7.76	-	-	-	-	-
9	10.34	12.76	83.42	89.00	6.53	-	-	-	-	-
10	76.26	72.22	6.49	24.24	70.97	0.04	0.24	94.97	99.99	9.65
15	-	-	-	-	-	0.02	0.80	94.68	99.96	9.58
Epoch	UGradSL					UGradSL+				
	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap (↓)	UA	MIA <sub>Score</sub>	RA	TA	Avg. Gap (↓)
10	14.98	33.22	77.18	84.07	13.27	6.26	14.10	93.39	99.62	4.94
11	24.26	34.38	68.22	75.06	20.37	6.52	11.66	93.04	99.37	5.51
12	28.70	24.62	68.17	74.39	19.22	21.46	27.38	89.41	97.07	9.85
13	38.46	72.90	61.78	64.72	37.75	29.48	31.92	87.74	94.93	12.88
14	99.86	86.74	0.45	0.20	88.02	31.62	32.68	86.53	93.36	13.51
Retrain	8.07	17.41	100.00	91.61	-	8.07	17.41	100.00	91.61	-