Label Smoothing Improves Machine Unlearning

Anonymous authors

000

001

003 004

010 011

012

013

014

015

016

017

018

019

021

024

025

026027028

029

031

033

034

036

038

040

041

042

043

044

046

047

048

051

052

Paper under double-blind review

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-andplay 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 also shows close connection to improve the local differential privacy. The consistent improvement in MU performance is only at a marginal cost of additional computations. For instance, UGradSL improves over the gradient ascent MU baseline constantly on different unlearning tasks without sacrificing unlearning efficiency. A self-adaptive UGradSL is also given for simple parameter selection.

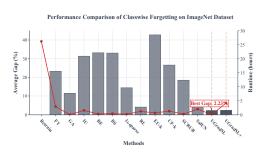
1 Introduction

Building a reliable ML model has become an important topic in this community. Machine unlearning (MU) is a task requiring to remove 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 product usage (Chen et al., 2021). 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 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 make the model unlearn the forgetting dataset without retraining the model.

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



055

057

064 065

066

067

068

069

071072073074

075

076

077

078 079

080

081

082

084

085

087

090

092

093

094

096

098

099

100

101

102

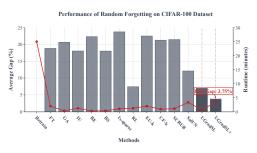
103

104

105

106

107



(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.

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 **Gra**dient-based **S**moothed **L**abels.

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; Bourtoule et al., 2021; Nguyen et al., 2022), gained prominence with the advent of privacy-focused legislation (Hoofnagle et al.) 2019; Pardau, 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.

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 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 θ_{tr} the original model trained on \mathcal{D}_{tr} and θ_u the model without the influence of \mathcal{D}_f . The goal of machine unlearning (MU) is how to generate θ_u from θ_{tr} .

Label Smoothing In a K-class classification task, let 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 $y_i^{\text{GLS},\alpha}$ with smooth rate $\alpha \in (-\infty,1]$ is $y_i^{\text{GLS},\alpha} = (1-\alpha) \cdot y_i + \frac{\alpha}{K} \cdot 1 = \left[\frac{\alpha}{K}, \cdots, \frac{\alpha}{K}, (1+\frac{1-K}{K}\alpha), \frac{\alpha}{K}, \cdots, \frac{\alpha}{K}\right]$, 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) θ_{tr}^* , the optimal parameters trained from $D_{tr} \sim \mathcal{D}_{tr}$, 2) θ_r^* , the optimal parameters trained from $D_r \sim \mathcal{D}_r$, such that $D_r = D_{tr} \setminus D_f$ and 3) θ_f^* , the optimal parameters unlearned using gradient ascent (GA) on $D_f \sim \mathcal{D}_f$. Note θ_r^* can be viewed as the *exact* MU model. The definitions of θ_{tr}^* and θ_r^* follow the standard empirical risk minimization as

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

and by using the influence function, θ_f^* is

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

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_{θ} is a function h parameterized by θ . ε is the weight of D_f compared with D_{tr} . The optimal parameter can be found when the gradient is 0:

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

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

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

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

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

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 are fixed. For GA, let $\varepsilon = -1$ in Eq. (3) and we have

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

where $(-\Delta\theta_r)$ represents the learning gap from θ_r^* to θ_{tr}^* while vector $\Delta\theta_f$ represents how much the model unlearns (backtracked progress) between θ_f^* and θ_{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 and the proof is given in Appendix C.1.

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

$$\sum_{z^f \in D_f} \nabla_{\boldsymbol{\theta}} \ell(h_{\boldsymbol{\theta}_r^*}, z^f) \approx -\boldsymbol{H}(\boldsymbol{\theta}_r^*, \boldsymbol{\theta}_{tr}^*) \cdot \sum_{z^f \in D_f} \nabla_{\boldsymbol{\theta}} \ell(h_{\boldsymbol{\theta}_{tr}^*}, z^f),$$

 $\boldsymbol{H}(\boldsymbol{\theta}_r^*, \boldsymbol{\theta}_{tr}^*) = \left[\sum_{z^{tr} \in D_{tr}} \nabla_{\boldsymbol{\theta}}^2 \ell(h_{\boldsymbol{\theta}_r^*}, z^{tr})\right] \left[\sum_{z^{r} \in D_r} \nabla_{\boldsymbol{\theta}}^2 \ell(h_{\boldsymbol{\theta}_{tr}^*}, z^{r})\right]^{-1}. \text{ Otherwise, there exist } \boldsymbol{\theta}_r^*, \boldsymbol{\theta}_{tr}^* \text{ such that } GA \text{ can not help MU, i.e., } \|\boldsymbol{\theta}_r^* - \boldsymbol{\theta}_f^*\| > \|\boldsymbol{\theta}_r^* - \boldsymbol{\theta}_{tr}^*\|.$

3.3 Label Smoothing Improves MU

Practically, we cannot guarantee that GA always helps MU as shown in Theorem 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

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

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

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

We leave the detailed form of $\Delta\theta_n$ to Eq. (21). But intuitively, $\Delta\theta_n$ captures the gradient influence of the smoothed non-target label on the weight. We show the effect of NLS (α < 0) in Theorem 2 below and its proof is given in Appendix (C.2).

Theorem 2. Given the approximation in Eq. (5) and $\langle \Delta \theta_r - \Delta \theta_f, \Delta \theta_n - \Delta \theta_f \rangle \le 0$, there exists an $\alpha < 0$ such that NLS improves GA in unlearning, i.e., $\|\theta_r^* - \theta_{f,NLS}^*\| < \|\theta_r^* - \theta_f^*\|$, where $\theta_{f,NLS}^*$ 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 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 4 in the Appendix.

The effect of the smoothed term in gradient ascent (GA) with NLS is equivalent to performing a gradient descent optimization with traditional 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 ϵ -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^{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 ϵ 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 ϵ -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 ϵ -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 α is more negative, the privacy of the labels in the forgetting dataset is better. When $\alpha \to (1 - \gamma_1/\gamma_2)$, we have $\epsilon \to 0$, indicating the best label-LDP result, which is the goal of MU. The theorem also warns that α cannot be arbitrarily negative.

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

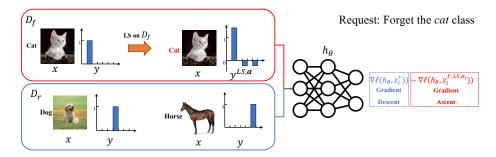


Figure 2: The framework of UGradSL. When there is an unlearning request, we can split the D_{tr} into D_f and D_r . We first apply label smoothing on $z_i^f = \{x,y\} \in D_f$ to get $z_i^{\text{LS},\alpha_i} = \{x,y^{\text{LS},\alpha_i}\}$, where the smooth rate can be pre-defined or self-adaptive.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

271

272

273

274

275

276

277

278

279

281

284

287

288

289

290291292

293

295

296

297

298

299

300

301

302

303

306 307

308

310

311

312

313

314

315

316

317318319

320 321

322

323

Algorithm 1 UGradSL+: A plug-and-play, efficient, gradient-based MU method using LS. UGradSL can be specified by imposing the dataset replacement in the bracket. If α is not given, the algorithm turns to the self-adaptive version.

Require: A almost-converged model $h_{\hat{\theta}_{tr}}$ trained with D_{tr} . The retained dataset D_r . The forgetting dataset D_f . Unlearning epochs E. GA ratio p. Distance threshold β . The optional smoothing ratio α .

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

How to choose the smooth rate α is worth discussion. Normally, the $\alpha_i \in \alpha$ for every data point $z_i^f \in D_f$ can be the same. To gain better performance, we improve UGradSL and UGradSL+ by taking every data point into consideration and assigning α_i individually and adaptively based on the distance $d(z_i^r, z_i^f) \in [0, 1]$ 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 lesser and should be reflected through a smaller α_i . The algorithm is presented in Algorithm and the framework is illustrated in Figure 2. We leave the details of the implementation, 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:

$$L(h_{\theta}, B_f^{\text{NLS}, \alpha}, B_r, p) = p \cdot \sum_{z^r \in B_r} \ell(h_{\theta}, z^r) - (1 - p) \cdot \sum_{\substack{z_i^f, \text{NLS}, \alpha_i \in B_f^{\text{NLS}, \alpha}}} \ell(h_{\theta}, z_i^f, \text{NLS}, \alpha_i)$$
(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. α is the vector for the smoothing rate of every data point z_i^f . h_{θ} 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 the LDP.

5 EXPERIMENTS AND RESULTS

5.1 EXPERIMENT SETUP

Dataset and Model Selection We validate our method using various datasets in different scales and modality, including CIFAR-10, CIFAR-100 (Krizhevsky et al., 2009), SVHN (Netzer et al., 2011),

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

	1		CIFA	AR-100			ImageNet							
Method	UA	MIAScore	RA	TA	Avg. Gap (↓)	RTE (\danh, min)	UA	MIA _{Score}	RA	TA	Avg. Gap ()	RTE (\dagger, hr)		
Retrain	100.00 _{±0.00}	$100.00_{\pm0.00}$	$99.96_{\pm0.01}$	$71.10_{\pm0.12}$	-	26.95	$100.00_{\pm0.00}$	$100.00_{\pm0.00}$	$71.62_{\pm 0.12}$	$69.57_{\pm0.07}$	-	26.18		
FT	$0.67_{\pm 0.38}$	$27.20_{\pm 1.34}$	$99.96_{\pm0.01}$	$71.46_{\pm0.09}$	43.12	1.74	$52.42_{\pm 15.81}$	$55.86_{\pm 18.02}$	$70.66_{\pm 2.54}$	$69.25_{\pm 0.78}$	23.25	2.87		
GA	$99.00_{\pm 0.57}$	$99.07_{\pm 0.50}$	$77.83_{\pm 2.07}$	$53.73_{\pm 0.96}$	10.36	0.06	$81.23_{\pm 0.69}$	$83.52_{\pm 2.08}$	$66.00_{\pm0.03}$	$64.72_{\pm 0.02}$	11.43	0.01		
IU	$2.07_{\pm 1.65}$	$33.20_{\pm 8.83}$	$99.96_{\pm0.01}$	$71.39_{\pm0.19}$	41.26	1.24	$33.54_{\pm 19.46}$	$49.83_{\pm 21.57}$	$66.25_{\pm 1.99}$	$66.28_{\pm 1.19}$	31.32	1.51		
BE	$99.07_{\pm 0.34}$	$99.00_{\pm0.49}$	$70.81_{\pm 2.69}$	$49.85_{\pm 1.32}$	13.08	0.55	$98.62_{\pm 0.58}$	$0.15_{\pm0.11}$	$53.13_{\pm0.27}$	$56.72_{\pm0.31}$	33.14	0.24		
BS	98.87 _{±0.57}	$98.73_{\pm 0.68}$	$71.16_{\pm 2.60}$	$50.03_{\pm 1.36}$	13.06	0.77	$98.85_{\pm0.50}$	$0.13_{\pm0.12}$	$53.35_{\pm0.16}$	$56.93_{\pm0.03}$	32.98	0.37		
ℓ_1 -sparse	$98.97_{\pm 1.03}$	$100.00_{\pm0.00}$	$86.99_{\pm0.76}$	$79.08_{\pm0.75}$	4.56	0.15	$100.00_{\pm0.00}$	$100.00_{\pm0.00}$	$39.01_{\pm 1.03}$	$44.62_{\pm 0.91}$	14.39	0.16		
RL	$99.80_{\pm 0.35}$	$100.00_{\pm 0.00}$	$99.97_{\pm 0.62}$	$77.31_{\pm 0.35}$	0.67	1.10	$100.00_{\pm 0.00}$	$100.00_{\pm 0.00}$	$62.06_{\pm 4.19}$	$62.93_{\pm0.45}$	4.05	1.17		
EU-k	$100.00_{\pm 0.00}$	$0.00_{\pm 0.00}$	$63.79_{\pm 1.10}$	$43.90_{\pm0.73}$	40.84	4.50	$100.00_{\pm 0.00}$	$0.00_{\pm 0.00}$	$32.99_{\pm 0.07}$	$37.19_{\pm0.15}$	42.75	0.62		
CF-k	$100.00_{\pm 0.00}$	$0.00_{\pm 0.00}$	94.88 ± 0.46	$61.32_{\pm 1.17}$	28.72	3.01	$99.79_{\pm 0.36}$	$0.00_{\pm 0.00}$	$66.84_{\pm0.03}$	$68.35_{\pm0.28}$	26.55	1.25		
SCRUB	$30.07_{\pm 49.48}$	$66.60_{\pm 29.19}$	$99.98_{\pm0.01}$	$77.97_{\pm 0.56}$	26.62	1.07	$56.59_{\pm 2.17}$	$75.59_{\pm 1.19}$	$66.98_{\pm0.11}$	$68.24_{\pm 0.07}$	18.45	0.21		
SalUN	$99.90_{\pm 0.01}$	$99.96_{\pm0.00}$	$99.98_{\pm 0.01}$	$75.02_{\pm0.10}$	1.02	2.15	$100.00_{\pm0.00}$	$100.00_{\pm0.00}$	$63.00_{\pm 5.03}$	$62.72_{\pm0.31}$	3.87	1.95		
UGradSL	$66.59_{\pm 0.90}$	$90.96_{\pm 5.05}$	$95.45_{\pm 1.42}$	$70.34_{\pm 1.78}$		0.07	$100.00_{\pm 0.00}$	$100.00_{\pm0.00}$	$76.91_{\pm 1.82}$			0.01		
UGradSL+	$100.00_{\pm 0.00}$	$100.00_{\pm0.00}$	$98.44_{\pm0.62}$	$74.12_{\pm 0.70}$	0.57	3.37	$100.00_{\pm 0.00}$	$100.00_{\pm 0.00}$	$78.16_{\pm0.07}$	$66.84_{\pm0.06}$	2.32	4.19		

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), ℓ_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 and CF-k (Goel et al., 2022). The implementation details of these baselines are given in Appendix E.1

Evaluation Metrics The evaluation metrics we use follows 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. 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 θ_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 MIA_{Score} = 1 - Pr $(x_f | \theta_*)$, where $x_f \in D_f$ are the forget samples and θ_* is the model under test. Overall, we use Avg. Gap to quantifies 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 3 respectively. The results of 20 NewsGroup and SVHN is given in Appendix 2.3 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, even in the randomness on 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.

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

			CII	FAR-100		Tiny-ImageNet						
Method	UA	MIAScore	RA	TA	Avg. Gap (↓)	RTE (↓, min)	UA	MIAScore	RA	TA	Avg. Gap (↓)	RTE (↓, min)
Retrain	$29.47_{\pm 1.59}$	$53.50_{\pm 1.19}$	$99.98_{\pm0.01}$	$70.51_{\pm 1.17}$	-	25.01	$49.35_{\pm0.38}$	$58.44_{\pm0.89}$	$83.80_{\pm0.29}$	$59.66_{\pm0.44}$	-	235.68
FT	$2.55_{\pm 0.03}$	$10.59_{\pm 0.27}$	$99.95_{\pm0.01}$	$75.95_{\pm0.05}$	18.83	1.95	29.23 _{±0.29}	$37.02_{\pm0.33}$	$82.51_{\pm0.20}$	$60.96_{\pm0.23}$	11.03	18.61
GA	$2.58_{\pm 0.06}$	$5.95_{\pm 0.17}$	$97.45_{\pm 0.02}$	$76.09_{\pm0.01}$	20.64	0.29	$19.34_{\pm 1.67}$	$25.19_{\pm 0.68}$	$81.51_{\pm 1.56}$	$59.66_{\pm0.61}$	16.39	8.65
IU	$15.71_{\pm 5.19}$	$18.69_{\pm 4.12}$	$84.65_{\pm 5.29}$	$62.20_{\pm 4.17}$	18.05	1.20	$60.61_{\pm 0.01}$	$83.67_{\pm0.15}$	$16.36_{\pm0.37}$	$23.44_{\pm0.29}$	35.04	7.30
BE	$0.01_{\pm 0.00}$	$1.45_{\pm 0.02}$	$99.97_{\pm0.18}$	$78.26_{\pm0.00}$	22.32	0.24	$17.65_{\pm0.31}$	$24.48_{\pm0.42}$	$82.85_{\pm0.20}$	$58.16_{\pm0.08}$	17.03	3.53
BS	$2.20_{\pm 1.21}$	$10.73_{\pm 9.37}$	$98.22_{\pm 1.26}$	$70.23_{\pm 1.67}$	18.02	0.34	$19.47_{\pm 0.69}$	$25.45_{\pm0.15}$	$81.23_{\pm 0.74}$	$56.75_{\pm0.80}$	17.09	5.63
ℓ_1 -sparse	$8.19_{\pm0.38}$	$19.11_{\pm 0.52}$	$88.39_{\pm0.31}$	$80.26_{\pm0.16}$	23.75	1.00	$35.73_{\pm 0.35}$	$41.98_{\pm0.73}$	$78.19_{\pm 0.05}$	$61.44_{\pm0.12}$	9.37	23.40
RL	$4.06_{\pm0.37}$	$50.12_{\pm 3.48}$	$99.92_{\pm 0.01}$	$71.30_{\pm0.12}$	7.41	1.20	$40.52_{\pm 0.15}$	$59.01_{\pm 0.76}$	$77.58_{\pm0.06}$	$60.18_{\pm0.19}$	4.04	27.08
EU-k	$1.73_{\pm 0.06}$	$3.33_{\pm 0.07}$	$98.44_{\pm 0.05}$	$59.92_{\pm0.43}$	22.51	1.96	$33.55_{\pm0.35}$	$22.19_{\pm 1.75}$	$81.41_{\pm 0.27}$	$58.08_{\pm0.21}$	14.01	20.02
CF-k	$0.07_{\pm 0.02}$	$0.47_{\pm 0.16}$	$99.98_{\pm0.01}$	$67.86_{\pm0.12}$	21.27	0.88	$19.31_{\pm 0.38}$	$23.22_{\pm 2.28}$	$81.59_{\pm0.37}$	$58.15_{\pm0.19}$	17.25	13.18
SCRUB	$0.09_{\pm 0.59}$	$4.01_{\pm 1.25}$	$99.97_{\pm 0.34}$	$77.45_{\pm0.26}$	21.46	1.06	$20.11_{\pm 1.15}$	$25.35_{\pm 7.53}$	$80.91_{\pm 0.77}$	$60.11_{\pm 0.99}$	16.42	25.79
SalUN	$35.23_{\pm0.32}$	$89.39_{\pm 0.46}$	$99.53_{\pm0.04}$	$64.26_{\pm 0.58}$	12.10	3.33	$40.39_{\pm0.15}$	$52.32_{\pm 10.67}$	$77.60_{\pm 0.11}$	$60.30_{\pm 0.31}$	5.48	34.42
UGradSL	18.36 _{±0.17}	$40.71_{\pm0.13}$	$98.38_{\pm0.03}$	$68.23_{\pm0.16}$	6.95	0.55	$40.73_{\pm0.71}$	$37.58_{\pm0.21}$	$67.30_{\pm0.04}$	$50.38_{\pm 0.77}$	13.82	9.47
UGradSL+	$21.69_{\pm 0.59}$	$49.47_{\pm 1.25}$	$99.87_{\pm 0.34}$	$73.60_{\pm0.26}$	3.75	3.52	$53.06_{\pm 1.27}$	$59.46_{\pm 1.01}$	$81.38_{\pm 0.75}$	$52.52_{\pm 0.84}$	3.57	25.93

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. 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.

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 bio-metric, 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.

5.3 Discussion

Influence Function in Deep Learning Influence function is 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 θ_r on CIFAR-10 dataset is given in Figure 5 in Appendix.

MIA as a Proxy for "Forgetfulness". Given a model θ_{\star} , 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 $A_f = \mathcal{A}(\theta_{\star}, \mathcal{D}_f)$ and $A_{te} = \mathcal{A}(\theta_{\star}, \mathcal{D}_{te})$, 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

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**.

			CI	FAR-20		I	l					
	UA	MIAScore	RA	TA	Avg. Gap (↓)	RTE (↓, min)	UA	MIAScore	RA	TA	Avg. Gap ()	RTE (↓, min)
Retrain	13.33 _{±1.64}	$28.47_{\pm 0.75}$	$99.94_{\pm0.01}$	$81.23_{\pm0.13}$	-	27.35	$6.74_{\pm 0.26}$	$9.77_{\pm 1.49}$	$94.38_{\pm0.49}$	$91.78_{\pm0.33}$	-	258.69
FT	$1.00_{\pm 0.43}$	$2.73_{\pm 0.52}$	99.37 _{±0.08}	$79.02_{\pm0.03}$	10.21	7.47	$5.36_{\pm0.17}$	$5.87_{\pm0.11}$	$93.91_{\pm 0.04}$	$93.18_{\pm0.03}$	1.79	25.94
GA	$87.93_{\pm 2.92}$	$88.93_{\pm 2.33}$	$81.46_{\pm 0.77}$	$64.07_{\pm 0.95}$	42.68	0.11	$6.00_{\pm0.16}$	$5.76_{\pm0.14}$	$92.86_{\pm0.13}$	$92.52_{\pm0.08}$	1.70	1.20
IU	$0.00_{\pm 0.00}$	$2.07_{\pm 1.29}$	$99.95_{\pm 0.01}$	$80.92_{\pm0.34}$	10.01	1.10	$5.90_{\pm 0.11}$	$4.91_{\pm 0.30}$	$93.05_{\pm0.01}$	$92.62_{\pm 0.01}$	1.97	219.77
BE	89.07 _{±1.39}	$91.73_{\pm 1.75}$	$76.36_{\pm0.92}$	$60.17_{\pm 0.92}$	45.91	0.33	$11.50_{\pm0.80}$	$48.41_{\pm 8.86}$	$88.37_{\pm0.81}$	$88.07_{\pm0.81}$	13.28	48.91
BS	$88.60_{\pm 1.13}$	$90.67_{\pm 1.18}$	$76.70_{\pm 1.08}$	$60.41_{\pm 1.17}$	45.38	0.29	$8.95_{\pm 5.11}$	$27.35_{\pm 30.20}$	$91.00_{\pm 5.22}$	$90.63_{\pm 5.65}$	6.08	50.99
ℓ_1 -sparse	$0.13_{\pm 0.09}$	$2.27_{\pm 0.57}$	$99.57_{\pm0.04}$	$80.44_{\pm0.08}$	10.14	0.38	$9.46_{\pm 1.82}$	$36.91_{\pm 30.96}$	$90.52_{\pm 1.75}$	$90.35_{\pm 1.77}$	8.79	37.49
RL	$56.93_{\pm 3.24}$	$98.60_{\pm0.29}$	$99.92_{\pm 0.01}$	$80.28_{\pm0.05}$	28.67	0.37	$8.31_{\pm0.43}$	$28.55_{\pm 16.74}$	$91.85_{\pm0.51}$	$91.62_{\pm0.42}$	5.76	40.09
EU-k	$8.00_{\pm 4.57}$	$16.33_{\pm 7.18}$	$97.07_{\pm0.18}$	$69.67_{\pm 0.35}$	7.98	0.87	$7.20_{\pm 0.19}$	$18.77_{\pm 3.69}$	$92.55_{\pm0.30}$	$91.04_{\pm 0.67}$	3.01	1.98
CF-k	$0.00_{\pm 0.00}$	$0.80_{\pm0.40}$	$99.98_{\pm0.01}$	$77.46_{\pm0.03}$	11.20	1.21	$5.46_{\pm0.32}$	$17.26_{\pm0.08}$	$94.45_{\pm0.04}$	$92.72_{\pm0.04}$	2.45	1.60
SCRUB	$0.00_{\pm 0.00}$	$1.13_{\pm0.34}$	$99.93_{\pm 0.01}$	$81.05_{\pm0.20}$	10.21	0.30	$8.78_{\pm0.77}$	$13.37_{\pm 5.22}$	$91.21_{\pm 0.86}$	$90.65_{\pm0.86}$	2.49	70.13
SalUN	$52.93_{\pm 2.21}$	$99.80_{\pm 0.35}$	$99.55_{\pm0.00}$	$76.48_{\pm0.26}$	29.02	2.88	$6.53_{\pm 0.28}$	$25.57_{\pm 8.22}$	$92.97_{\pm 0.03}$	$92.27_{\pm 0.07}$	4.48	83.43
UGradSL	22.87 _{±0.90}	$38.93_{\pm 1.57}$	$97.20_{\pm0.19}$	$75.84_{\pm0.16}$	7.03	0.13	$6.29_{\pm 1.41}$	$5.73_{\pm 3.50}$	$93.44_{\pm0.14}$	$92.80_{\pm0.27}$	1.61	2.17
UGradSL+	$78.44_{\pm 1.19}$	$88.67_{\pm 0.35}$	$97.93_{\pm0.71}$	$79.77_{\pm 0.58}$	32.20	8.12	$6.12_{\pm0.31}$	$5.54_{\pm 0.34}$	$92.79_{\pm 0.01}$	$92.49_{\pm 0.04}$	1.79	51.41

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 the same. The plots of loss distribution for random and class-wise forgetting are given in Figure 6 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.8.

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	MIAScore	RA	TA	Add. MIA	Avg. Gap (\$)	RTE (\dagger, min)	UA	MIAScore	RA	TA	Add. MIA	Avg. Gap (\$)	RTE (↓, min)		
Retrain	$100.00_{\pm0.00}$	$100.00_{\pm0.00}$	$98.19_{\pm 3.14}$	$94.50_{\pm0.34}$	$99.23_{\pm 0.08}$	-	24.62	$8.07_{\pm0.47}$	$17.41_{\pm 0.69}$	$100.00_{\pm0.01}$	$91.61_{\pm0.24}$	$50.69_{\pm 0.73}$	-	24.66		
FT	$22.71_{\pm 5.31}$	$79.21_{\pm 8.60}$	$99.82_{\pm0.09}$	$94.13_{\pm0.14}$	99.09 _{±0.07}	20.04	2.02	$1.10_{\pm0.19}$	$4.06_{\pm0.41}$	$99.83_{\pm0.03}$	$93.70_{\pm0.10}$	$54.05_{\pm0.31}$	5.19	1.58		
GA	$25.19_{\pm 11.38}$	$73.48_{\pm 9.68}$	$96.84_{\pm0.58}$	$73.10_{\pm 1.62}$	$99.43_{\pm 0.09}$	24.86	0.08	$0.56_{\pm0.01}$	$1.19_{\pm 0.05}$	$99.48_{\pm0.02}$	$94.55_{\pm 0.05}$	$55.04_{\pm0.66}$	6.31	0.31		
IU	$83.92_{\pm 1.16}$	$92.59_{\pm 1.41}$	$98.77_{\pm0.12}$	$92.64_{\pm0.23}$	$99.71_{\pm 0.07}$	5.28	1.18	$17.51_{\pm 2.19}$	$21.39_{\pm 1.70}$	$83.28_{\pm 2.44}$	$78.13_{\pm 2.85}$	$53.98_{\pm 0.55}$	9.37	1.18		
BE	$64.93_{\pm0.01}$	$98.19_{\pm0.00}$	$99.47_{\pm 0.00}$	$94.00_{\pm0.11}$	$99.60_{\pm 0.02}$	7.81	0.20	$0.00_{\pm 0.00}$	$0.26_{\pm0.02}$	$100.00_{\pm0.00}$	$95.35_{\pm0.18}$	$55.41_{\pm 0.49}$	6.74	3.17		
BS	$93.69_{\pm 4.32}$	$99.82_{\pm0.04}$	$97.69_{\pm 1.29}$	$92.89_{\pm 1.26}$	$99.56_{\pm0.10}$	1.79	0.29	$0.48_{\pm 0.07}$	$1.16_{\pm0.04}$	$99.47_{\pm0.01}$	$94.58_{\pm0.03}$	$55.88_{\pm0.72}$	6.51	1.41		
ℓ_1 -sparse	$100.00_{\pm0.00}$	$100.00_{\pm0.00}$	$97.86_{\pm 1.29}$	$96.11_{\pm 1.26}$	$99.02_{\pm 0.15}$	0.43	1.00	$2.80_{\pm0.37}$	$18.59_{\pm 3.48}$	$99.97_{\pm0.01}$	$94.08_{\pm0.12}$	$52.17_{\pm 0.87}$	2.09	1.98		
RL	$99.99_{\pm 0.01}$	$100.00_{\pm0.00}$	$100.00_{\pm0.00}$	$95.50_{\pm0.11}$	$99.08_{\pm 0.07}$	0.59	1.04	$2.80_{\pm0.37}$	$18.59_{\pm 3.48}$	$99.97_{\pm0.01}$	$94.08_{\pm0.12}$	$52.17_{\pm 0.87}$	2.09	1.98		
EU-k	$100.00_{\pm0.00}$	$100.00_{\pm0.00}$	$100.00_{\pm0.00}$	$75.04_{\pm 1.10}$	$99.89_{\pm 0.18}$	4.39	1.45	$0.00_{\pm 0.00}$	$0.50_{\pm0.30}$	$99.99_{\pm0.01}$	$77.21_{\pm 1.21}$	$61.88_{\pm 1.33}$	10.12	1.58		
CF-k	$100.00_{\pm0.00}$	$100.00_{\pm0.00}$	$100.00_{\pm0.00}$	$78.95_{\pm0.53}$	$100.00_{\pm 0.00}$	3.63	1.32	$0.00_{\pm 0.00}$	$0.00_{\pm 0.00}$	$100.00_{\pm0.00}$	$80.98_{\pm0.27}$	69.91+1.33	11.07	1.47		
SCRUB	$100.00_{\pm0.00}$	$100.00_{\pm0.00}$	$99.93_{\pm0.01}$	$95.22_{\pm 0.07}$	$100.00_{\pm0.00}$	0.65	1.09	$0.70_{\pm 0.59}$	$3.88_{\pm 1.25}$	$99.59_{\pm0.34}$	$94.22_{\pm 0.26}$	55.33+0.59	5.71	4.05		
SalUN	$90.74_{\pm 13.91}$	$100.00_{\pm0.00}$	$98.20_{\pm0.34}$	$80.49_{\pm 1.21}$	$98.63_{\pm0.59}$	4.78	2.22	$46.95_{\pm0.15}$		$97.75_{\pm0.42}$		$69.95_{\pm0.12}$	28.74	2.42		
UGradSL	94.99 _{±4.35}	$97.95_{\pm 1.78}$	95.47 _{±4.08}	86.78 _{±5.68}	$99.94_{\pm0.01}$	3.64	0.22	$5.87_{\pm 0.51}$	$13.33_{\pm 0.70}$	$98.82_{\pm0.28}$		$53.54_{\pm 0.97}$	1.78	0.45		
UGradSL+	$100.00_{\pm0.00}$	$100.00_{\pm0.00}$	$99.26_{\pm0.01}$	$94.29_{\pm 0.07}$	$100.00_{\pm 0.00}$	0.41	3.07	$6.03_{\pm0.17}$	$10.65_{\pm0.13}$	$99.79_{\pm 0.03}$	$93.64_{\pm0.16}$	$52.29_{\pm 0.85}$	2.53	3.07		

The study of the Streisand effect Jansen & Martin (2015) and gradient analysis are given in Appendix E.6 and E.7 respectively.

6 CONCLUSIONS AND LIMITATIONS

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. Our work has limitations. For example, we desire an efficient way to find the exact MU state in experiments and further explore the applications of MU to promote privacy and fairness. Our method can be further validated and tested in other tasks, such as unlearning recommendation systems, etc.

REFERENCES

- Martín Abadi, Andy Chu, Ian Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. *arXiv: Machine Learning*, 2016. doi: 10.1145/2976749.2978318.
 - 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 Bourtoule, 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 Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine J. Lee, Katherine Lee, Adam Roberts, Tom B. Brown, Dawn Song, Úlfar Erlingsson, Úlfar Erlingsson, Úlfar Erlingsson, Ulfar Erlingsson, Alina Oprea, and Colin Raffel. Extracting training data from large language models. *arXiv: Cryptography and Security*, 2020. doi: null.
 - Min Chen, Zhikun Zhang, Zhikun Zhang, Tianhao Wang, Tianhao Wang, Michael Backes, Michael Backes, Mathias Humbert, Yang Zhang, and Yang Zhang. When machine unlearning jeopardizes privacy. *arXiv: Cryptography and Security*, 2020. doi: 10.1145/3460120.3484756.
 - 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, 2021.
 - 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.
 - Christopher A. Choquette-Choo, Christopher A. Choquette Choo, Christopher A. Choquette Choo, Christopher A. Choquette-Choo, Florian Tramèr, Florian Tramèr, Nicholas Carlini, Nicolas Papernot, Nicolas Papernot, and Nicolas Papernot. Label-only membership inference attacks. *International Conference on Machine Learning*, 2020. doi: null.
 - 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, Matt Fredrikson, Somesh Jha, Somesh Jha, Thomas Ristenpart, and Thomas Ristenpart. Model inversion attacks that exploit confidence information and basic countermeasures. *Conference on Computer and Communications Security*, 2015. doi: 10.1145/2810103.2813677.
 - Karan Ganju, Karan Ganju, Qi Wang, Qi Wang, Wei Yang, Wei Yang, Carl A. Gunter, Carl A. Gunter, Nikita Borisov, and Nikita Borisov. Property inference attacks on fully connected neural networks using permutation invariant representations. *Conference on Computer and Communications Security*, 2018. doi: 10.1145/3243734.3243834.
 - 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, Sorami Hisamoto, Matt Post, Matt Post, Kevin Duh, and Kevin Duh. Membership inference attacks on sequence-to-sequence models. *arXiv: Learning*, 2019. doi: 10.1162/tacl_a_00299.

- Nils Homer, Nils Homer, Szabolcs Szelinger, Szabolcs Szelinger, Margot Redman, Margot Redman, David Duggan, David Duggan, Waibhav Tembe, Waibhav Tembe, Jill Muehling, Jill Muehling, John V. Pearson, John V. Pearson, Dietrich A. Stephan, Dietrich A. Stephan, Stanley F. Nelson, Stanley F. Nelson, David W. Craig, and David Craig. Resolving individuals contributing trace amounts of dna to highly complex mixtures using high-density snp genotyping microarrays. *PLOS Genetics*, 2008. doi: 10.1371/journal.pgen.1000167.
- 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.
- 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.
- K. Rustan M. Leino, Klas Leino, Matt Fredrikson, and Matt Fredrikson. Stolen memories: Leveraging model memorization for calibrated white-box membership inference. *USENIX Security Symposium*, 2019. doi: null.
- Zheng Li, Zheng Li, Yang Zhang, and Yang Zhang. Membership leakage in label-only exposures. *Conference on Computer and Communications Security*, 2021. doi: 10.1145/3460120.3484575.
- 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.

- Michal Lukasik, 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, H. Brendan McMahan, Daniel Ramage, Daniel Ramage, Kunal Talwar, Kunal Talwar, Li Zhang, and Li Zhang. Learning differentially private recurrent language models. *International Conference on Learning Representations*, 2018. doi: null.
 - 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. *IEEE Symposium on Security and Privacy*, 2019. doi: 10.1109/sp.2019.00065.
 - 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 Pardau. 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, Md. Atiqur Rahman, Md. A. Rahman, Md. Atiqur Rahman, Tanzila Rahman, Tanzila Rahman, Robert Laganière, Robert Laganière, Robert Laganière, Noman Mohammed, and Noman Mohammed. Membership inference attack against differentially private deep learning model. *Transactions on Data Privacy*, 2018. doi: null.
 - Alexandre Sablayrolles, Alexandre Sablayrolles, Matthijs Douze, Matthijs Douze, Yann Ollivier, Yann Ollivier, Cordelia Schmid, Cordelia Schmid, Hervé Jegou, and Hervé Jegou. White-box vs black-box: Bayes optimal strategies for membership inference. *arXiv: Machine Learning*, 2019. doi: null.
 - 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, Reza Shokri, Vitaly Shmatikov, and Vitaly Shmatikov. Privacy-preserving deep learning. *Allerton Conference on Communication, Control, and Computing*, 2015. doi: 10.1145/2810103.2813687.
 - Reza Shokri, Reza Shokri, Marco Stronati, Marco Stronati, Congzheng Song, Congzheng Song, Vitaly Shmatikov, and Vitaly Shmatikov. Membership inference attacks against machine learning models. *IEEE Symposium on Security and Privacy*, 2017. doi: 10.1109/sp.2017.41.

- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
 - Liwei Song, Liwei Song, Prateek Mittal, and Prateek Mittal. Systematic evaluation of privacy risks of machine learning models. *arXiv: Cryptography and Security*, 2020. doi: null.
 - Thomas Steinke, Milad Nasr, and Matthew Jagielski. Privacy auditing with one (1) training run. *Neural Information Processing Systems*, 2023. doi: 10.48550/arxiv.2305.08846.
 - 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, Samuel Yeom, Irene Giacomelli, Irene Giacomelli, Matt Fredrikson, Matt Fredrikson, Somesh Jha, Somesh Jha, and Somesh Jha. Privacy risk in machine learning: Analyzing the connection to overfitting. *IEEE Computer Security Foundations Symposium*, 2018. doi: 10.1109/csf.2018.00027.