

Towards Safe Machine Unlearning: a Paradigm that Mitigates Performance Degradation

Anonymous Author(s)

Abstract

We study the machine unlearning problem which aims to remove specific training data from a pre-trained machine learning model to allow users to exercise their ‘right to be forgotten’ to protect user privacy. Conventional machine unlearning methods would degrade the model performance after the unlearning procedure. To mitigate the issue, they typically rely on the access to the remaining training data to fine-tune the unlearned model to mitigate the influence of unlearning. However, accessing the remaining training data may not always be practical for different reasons (e.g., data expiration policies, storage limitations, or additional privacy constraints). Machine unlearning without access to the remaining training data poses significant challenges to retaining model performance. In this paper, we study how to unlearn specific training data from a pre-trained model without accessing the remaining training data and protect model performance without dramatically changing the model’s parameters. We propose a practical method called Targeted Label Noise Injection. Intuitively, our method assigns incorrect yet controllable labels to the examples that need to be forgotten and fine-tunes the pre-trained model to learn these new labels. This strategy effectively moves the to-be-forgotten examples across the decision boundary with a small impact on the model’s overall performance. We theoretically prove the effectiveness of the proposed method and empirically show that it achieves state-of-the-art unlearning performance across various datasets.

Relevance: Machine unlearning is critical to user modeling because users have the ‘right to be forgotten’. This paper proposes a novel paradigm for machine unlearning to mitigate model performance degradation during unlearning, which is highly relevant to the track of ‘User modeling, personalization and recommendation’. Moreover, the proposed method provides new state-of-the-art for Web applications where both high privacy and utility are required.

CCS Concepts

- Networks → Network privacy and anonymity; • Security and privacy → Usability in security and privacy.

Keywords

Privacy protection, machine unlearning, trustworthy machine learning

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

WWW’25, April 28th–May 2nd 2025, Sydney, NSW, Australia

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-XXXX-X/18/06

<https://doi.org/XXXXXXX.XXXXXXX>

ACM Reference Format:

Anonymous Author(s). 2018. Towards Safe Machine Unlearning: a Paradigm that Mitigates Performance Degradation. In *Proceedings of the ACM Web Conference 2025 (WWW’25), April 28th–May 2nd, 2025, Sydney, NSW, Australia*. ACM, New York, NY, USA, 23 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 Introduction

In the era of big data, machine learning models have become integral to a wide range of web applications, from personalized recommendations to critical decision-making systems. These models are often trained on vast amounts of user data, much of which may be sensitive or personal. As concerns about data privacy escalate, stringent regulations such as the General Data Protection Regulation (GDPR) [27] have been enacted, granting individuals the right to have their personal data erased/forgotten. There is a growing need for machine learning models to not only learn from data but also to unlearn specific data points upon request.

Machine unlearning [2] aiming to remove specific data from trained machine learning models has thus emerged as a crucial capability. This ensures compliance with data privacy laws and maintains user trust in systems that handle personal information. Unfortunately, the unlearning procedure would degrade model performance. To address this issue, conventional methods typically require access to both the data to be forgotten and the remaining training data, using the remaining data to tune the unlearned model to mitigate performance degradation [1, 5, 35]. However, accessing the remaining data can be impractical due to many reasons, e.g., data expiration policies that mandate deletion after a certain period, storage limitations that prevent retaining large datasets, or additional privacy constraints that prohibit the reuse of data without explicit consent.

Unlearning without access to the remaining data poses significant challenges. Existing machine unlearning methods in this setting (e.g., [4, 8, 23]) are limited and often lead to substantial changes in the model’s parameters, which can degrade performance on the remaining data that should be retained. Intuitively, fine-tuning solely on the examples to be forgotten may require extensive parameter adjustments, resulting in overfitting the to-be-forgotten examples and underperformance on the remaining examples. Motivated by these challenges, in this paper, we study how to safely unlearn some examples from a pre-trained model without accessing the remaining training examples and without dramatically changing the model’s parameters.

To address this problem, we propose a novel method called Targeted Label Noise Injection for machine unlearning. The key philosophy of our method is to minimize parameter changes by selecting incorrect yet easy-to-learn labels for the to-be-forgotten examples. Specifically, we assign each instance to be forgotten an incorrect label that the model is already somewhat inclined toward. This ensures that the pseudo-label aligns relatively well with the model’s existing knowledge and is easy for the model to learn. By doing

so, we reduce the need for large parameter updates, helping to maintain the model's performance on the remaining examples.

However, it is also important to note that during fine-tuning, the model can overfit the pseudo-labels with unnecessarily high confidence. This leads to unwanted parameter adjustments. To prevent overfitting the pseudo-labels, we reduce the contribution of each to-be-forgotten example to the loss once the model starts predicting its pseudo-label correctly. This is achieved by introducing an exponent λ to the training loss. By setting a high value for λ during the fine-tuning process, the loss of instances with pseudo-labels, which the model has already learned, is dramatically reduced. Consequently, in the later fine-tuning epochs, these instances have little influence on the model's parameter updates, effectively avoiding unnecessary changes to the model's parameters.

Furthermore, we theoretically analyze the parameter changes induced by our method and find that it concentrates updates on a specific subset of parameters associated with to-be-forgotten examples. This targeted adjustment effectively modifies only the necessary parameters while minimizing the impact on unrelated parts of the model. Empirically, we have validated that our method achieves the smallest parameter changes across different datasets when compared with various baselines. Our method consistently achieves state-of-the-art performance for machine unlearning without accessing the remaining data, demonstrating its effectiveness and practicality across diverse scenarios.

2 Related Work

We study the machine unlearning aiming to enable a pre-trained model to discard information acquired from a subset of data. The primary objective has been to undo the influence of undesired data on the model while maintaining the model's predictive power on the remaining sample. Historically, simpler machine learning models, such as linear/logistic regression, k-means clustering, and random forests, have been the subjects of various unlearning techniques (e.g., [10, 24, 25]). However, these methodologies are inherently designed for these less complex models and don't seamlessly translate to intricate architectures like deep neural networks.

To achieve unlearning for deep neural networks, many methods have also been proposed. Some of them focus on making the model forget all examples corresponding to a specific class, like all images of a particular category in a dataset [5, 35, 36]. Some of them emphasize erasing information from some examples potentially spanning multiple classes from the pre-trained model [3, 26, 29]. By considering that data deletion may requests are practically received on a per-instance basis, some methods which allow to forget any subset of training examples have also been proposed [9, 13, 37]. A notable limitation shared by most method is their reliance on access to both the data intended for unlearning and insights on the remaining training data during the unlearning phase (e.g., [17, 21, 31]). This dependency becomes problematic in real-world scenarios where accessing the remaining sample might be hindered by factors like data expiration policies or other regularization.

To circumvent this challenge, the use of gradient ascent in the unlearning process has been naturally proposed. The basic idea is to employ gradient ascent to reverse the loss value of an example marked for unlearning from a pre-trained model. However, while

the gradient ascent-based unlearning method has demonstrated its effectiveness across diverse settings and datasets [16, 28, 33], concerns have been raised about its suitability in scenarios demanding high privacy and utility. Graves et al. [12] argued that merely amplifying loss using gradient ascent doesn't inherently ensure heightened privacy. Suriyakumar and Wilson [30] further highlighted the potential for data leakage even when the loss is accentuated. This brings to light the privacy implications of deploying existing techniques. Some methods have also been proposed which use gradient descent (e.g., [16, 22, 32]). However, fine-tuning solely on the examples to be forgotten may require extensive parameter adjustments. This leads to overfitting to the to-be-forgotten examples and underperformance on remaining examples. Motivated by these challenges, we study how to safely unlearn examples from a pre-trained model without accessing the remaining training examples and without dramatically changing the model's parameters.

3 δ -Targeted^λ Label Noise Injection for Machine Unlearning

Problem Setup. Consider a dataset $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, where $\mathbf{x}_i \in \mathcal{X}$ is an instance from the feature space, $y_i \in \{1, 2, \dots, C\}$ is its corresponding label from C classes, n is the size of the dataset. Each example (\mathbf{x}_i, y_i) is assumed to be independently drawn from an underlying distribution $P(X, Y)$.

Let $\sigma \circ f_{\hat{\theta}} : \mathcal{X} \rightarrow \mathbb{R}^C$ be a pre-trained model trained on D by minimizing the average loss, where $f_{\hat{\theta}}$ represents the deep network and $\hat{\theta}$ denotes the trained parameters, and σ is the softmax function. Given an instance \mathbf{x} , the model outputs $\sigma \circ f_{\hat{\theta}}(\mathbf{x})$ estimates the class-posterior distribution.

Suppose we have a subset $D_f \subseteq D$ of examples that we aim to forget from the model. The remaining examples are denoted as $D_r = D \setminus D_f$. Our objective is to obtain a model that effectively forgets D_f while maintaining high accuracy on D_r and generalization, using only D_f and the pre-trained model $\sigma \circ f_{\hat{\theta}}$.

3.1 Targeted Label Noise Injection with Probability δ

Our method for machine unlearning is based on injecting label noise into the to-be-forgotten examples. Intuitively, to mitigate performance degradation, we select incorrect yet easy-to-learn labels for the to-be-forgotten examples. This approach ensures that the pseudo-label aligns relatively well with the model's existing knowledge and is easy to learn. This effectively reduces the need for large parameter changes and helps to maintain the model's performance.

Specifically, for each to-be-forgotten example $(\mathbf{x}, y) \in D_f$, we compute the model's predicted probabilities $\sigma \circ f_{\hat{\theta}}(\mathbf{x})$. Let $\hat{y} = \arg \max_{i \in \{1, 2, \dots, C\}} \sigma \circ f_{\hat{\theta}}(\mathbf{x})$ be the predicted label. The second most confident label is defined by $\hat{y}' = \arg \max_{i \in \{1, 2, \dots, C\} \setminus \hat{y}} \sigma \circ f_{\hat{\theta}}(\mathbf{x})$. It is easy to find that the second most confident label of a to-be-forgotten example is an easy-to-learn label for it. In this paper, we use the second most confident labels as the easy-to-learn labels.

By considering that different scenarios may require different degrees of forgetting, we might want to control the trade-off between

the level of unlearning and the preservation of the models's performance. Additionally, if all to-be-forgotten examples are always assigned incorrect labels, an adversary could infer that the true labels are different from the predicted ones, thus gaining information about the original data. Therefore, we introduce a parameter $\delta \in [0, 1]$ to control the strength of the label perturbation. Specifically, if the model's predicted label \hat{y} does not match the original label y , this indicates a misprediction. In such a case, with a probability of $1 - \delta$, the pseudo label y' is set to the predicted label \hat{y} ; and with a probability of δ , the original label y is retained as the pseudo label y' . Conversely, if the predicted label \hat{y} matches the original label y , indicating a correct prediction, then with a probability of $1 - \delta$, the pseudo label y' is set to the second most confident label \hat{y}' ; and with a probability of δ , the original label y is retained as the pseudo label y' . Formally, this assignment of the pseudo label y' can be expressed as:

$$y' := \begin{cases} y, & \text{if } \alpha < \delta; \\ \hat{y}', & \text{if } \hat{y} = y \wedge \alpha \geq \delta; \\ \hat{y}, & \text{if } \hat{y} \neq y \wedge \alpha \geq \delta, \end{cases}$$

where $\alpha \sim \text{Uniform}(0, 1)$ is a random variable. By choosing a smaller δ , we enforce that more pseudo-labels differ from the original labels. To let the pre-trained model randomly guess a to-be-forgotten example, δ should be set to $1/C$, where C represents the number of unique labels.

3.2 Avoid Overfitting to Label Noise with a Loss Exponent λ

After injecting label noise into the to-be-forgotten examples. A set of new to-be-forgotten examples $D'_f = \{(x_i, y'_i)\}_{i=1}^m$ are obtained which contains instances and their pseudo labels¹, where m is the size of example to be forgotten. The unlearning process aims to minimize the loss induced by the pseudo labels while preserving the accuracy on the remaining sample D_r by only the sample D_f . Note that since the pseudo labels are mostly different from the original labels, small losses with respect to the pseudo labels imply large losses with respect to the original labels. Examples with large loss values mean that they have not been fitted by the model and are thus unlearned.

However, we notice that minimizing the loss natively might lead the model to fit the pseudo labels with unnecessarily high confidence, which is not our primary objective. Our goal is to make the model mispredict the original labels of these instances using pseudo labels, but without letting the model overfit these pseudo labels with high confidence. This distinction is crucial. If the model already predicts a pseudo label for a to-be-forgotten example correctly, there's no need to adjust its parameters further to increase the confidence of this prediction. Instead, during the unlearning procedure, the model updating should focus on the pseudo labels that it cannot predict correctly. By doing so, we can avoid unnecessary parameter adjustments, which could lead to overfitting the to-be-forgotten examples and adversely affect the model's performance on remaining examples.

¹Note that in the rest of the paper when there is no confusion, we do not distinguish to-be-forgotten examples with original labels and pseudo label.

Algorithm 1 δ -Targeted γ Label Noise Injection For Machine Unlearning

```

1: procedure UNLEARN( $D_f, \hat{\theta}, \eta, \delta$ )
2:    $D'_f \leftarrow \emptyset$                                  $\triangleright$  Initialize an empty set
3:   for each  $(x, y) \in D_f$  do  $\triangleright$  Iterate over all examples in the
   to-be-forgotten examples
4:      $\hat{y} \leftarrow \arg \max_i \sigma \circ f_{\hat{\theta}}(x)$   $\triangleright$  The model's predicted label
5:      $\hat{y}' \leftarrow \arg \max_{i \neq \hat{y}} \sigma \circ f_{\hat{\theta}}(x)$        $\triangleright$  The model's second
   confident label
6:      $\alpha \sim \mathcal{U}(0, 1)$                                  $\triangleright$  Draw a scalar from the uniform
   distribution.
7:     if  $\hat{y} \neq y \wedge \alpha < \delta$  then
8:        $y' \leftarrow y$                                       $\triangleright$  Retain the original label
9:     else if  $\hat{y} \neq y \wedge \alpha \geq \delta$  then
10:       $y' \leftarrow \hat{y}$                                      $\triangleright$  Let the pseudo label be the model's
    predicted label.
11:    else if  $\hat{y} = y \wedge \alpha < \delta$  then
12:       $y' \leftarrow y$                                       $\triangleright$  Retain the original label
13:    else
14:       $y' \leftarrow \hat{y}'$                                  $\triangleright$  Let the pseudo label be the model's
    second confident label.
15:    end if
16:     $D'_f \leftarrow D'_f \cup \{(x, y')\}$   $\triangleright$  Add new example  $(x, y')$  to  $D'_f$ 
17:  end for
18:   $\hat{\theta}'' \leftarrow \hat{\theta}$   $\triangleright$  Initialize the parameter needed to be fine-tuned
19:  while  $L'(\hat{\theta}'', D'_f) \geq \epsilon$  do
20:    Compute  $\nabla L'(\hat{\theta}'', D'_f)$  using Eq. (1)           $\triangleright$  Compute
   gradient of modified loss
21:     $\hat{\theta}'' \leftarrow \hat{\theta}'' - \eta \nabla L'(\hat{\theta}'', D'_f)$        $\triangleright$  Update the model's
   parameters
22:  end while
23:  return  $\hat{\theta}''$                                       $\triangleright$  Return the updated model's parameters
24: end procedure

```

A natural approach is to dynamically remove examples from the training set once their pseudo labels can be accurately predicted by a model during the unlearning process. This ensures that these examples are not further learned. The model should then concentrate on learning the pseudo labels of the remaining to-be-forgotten examples that it cannot predict accurately. However, we found that this should not work well in practice. The primary issue is that learning the pseudo labels for the remaining examples alters the model's parameters. After these changes, the model may not remember the pseudo labels of previously removed examples from the to-be-forgotten examples. Therefore, instead of removing an example from the training set once its pseudo label can be accurately predicted, a more effective method is to reduce the example's contribution to the optimizing objective if the example already exhibits a small loss.

Motivated by this, we propose using an exponent γ to control the loss contribution of to-be-forgotten examples during the learning process. If the model already aligns well with the pseudo label of a to-be-forgotten example, instead of removing this example from the training set, we reduce its contribution to the loss with an exponent γ , i.e., ℓ^γ . Note that if $\gamma > 1$, when ℓ is close to zero, ℓ^γ will become

smaller than ℓ ; while when ℓ is larger than one, ℓ^y will become larger than ℓ . If the model predicts the pseudo label of a to-be-forgotten example, the corresponding loss value should be close to zero. If the model struggles to predict the pseudo label of a to-be-forgotten example, that example will have a larger loss compared with other examples whose labels are correctly predicted, encouraging the model to focus more on learning such an unfitted example. This method automatically adjusts the influence of each to-be-forgotten example based on its prediction accuracy. By decreasing the loss contribution from well-predicted examples and increasing it for poorly predicted ones during the unlearning process, this method ensures more efficient and targeted unlearning.

Formally, we exponentiate cross-entropy loss function $\ell(\sigma \circ f_\theta(\mathbf{x}), \mathbf{y}')$ by a exponent λ , where \mathbf{y}' is the one-hot encoding form of pseudo label y' . This implicitly introduces a dynamic weight associated with each example. The objective function is defined as:

$$\arg \min_{\theta} L'(\theta, D'_f) = \arg \min_{\theta} \frac{1}{m} \sum_{(\mathbf{x}, \mathbf{y}') \in D'_f} \frac{\ell(\sigma(f_\theta)(\mathbf{x}), \mathbf{y}')^\lambda}{m}. \quad (1)$$

By raising the loss to the power of λ , we dynamically adjust the contribution of each instance to the overall objective. This method ensures that the model concentrates on learning the pseudo-labels it hasn't yet mastered, without over-adjusting parameters for those it already predicts correctly. By preventing overfitting to the pseudo-labels, we avoid unnecessary parameter changes that could degrade performance on the remaining data. The gradient update during fine-tuning is then

$$\theta \leftarrow \theta - \eta \nabla L'(\theta, D'_f), \quad (2)$$

where η is the learning rate. By combining the above with the targeted label noise injection, our unlearning strategy effectively reduces the risk of over-adjusting the model parameters during the fine-tuning process. We name our method δ -Targeted ^{λ} Sample Unlearning, the pseudo-code is illustrated in Algorithm 1.

4 Theoretical Analysis

In this section, we provide a theoretical analysis demonstrating how our method concentrates parameter updates on the specific subset of parameters associated with the to-be-forgotten examples. This targeted adjustment effectively modifies the necessary parameters while reducing the impact on unrelated parts of the model. We compare our method with sample unlearning by applying gradient ascent on original labels. Note that gradient ascent is commonly employed by existing methods (e.g., [11, 22, 32]). In experiments, we show that benefiting from the concentrated parameter updates, our method leads to smaller changes in parameters compared with all baselines on different neural network structures. As a consequence, the smaller change in parameter further lead to the small performance degradation

We formally introduce Concentrated Parameter Updates in the following definition.

DEFINITION 4.1 (CONCENTRATED PARAMETER UPDATES). Let θ denote the original parameter set of a pre-trained model, and let θ_1 and θ_2 be two distinct sets of updated parameters. Define the changes in parameters for θ_1 and θ_2 relative to θ as $\Delta\theta_1 = \theta_1 - \theta$ and $\Delta\theta_2 = \theta_2 - \theta$, respectively. We say that θ_2 represents a more concentrated parameter update than θ_1 if

- (1) $\|\Delta\theta_1\|_1 = \|\Delta\theta_2\|_1$, and
- (2) $\Delta\theta_2$ exhibits larger changes for certain parameters and smaller influence for others, compared to a more uniformly distributed change in $\Delta\theta_1$.

Note that the above definition restricts the total magnitude of change $\|\Delta\theta_1\|_1 = \|\Delta\theta_2\|_1$ across two sets of parameters. It ensures that comparisons are not biased by the scale of the changes, thereby facilitating a fair comparison of parameter updates based on their distribution of changes.

Given the non-convex optimization landscape of neural networks and their inherent non-linearity, a fully transparent statistical analysis comparing the effects of gradient ascent and descent to all parameters is ambitious. Empirically, we have found that changes in the last fully connected layer serve as an effective surrogate for monitoring changes in the model's parameters. Specifically, let ϕ denote ϕ as the parameter of the last fully connected layer, in Section 5.1, we demonstrate that variations in ϕ exhibit a strong dependence with changes spanning all parameters². Hence, we narrow our analysis to examine the impact of the unlearning process on the last fully connected layer's parameters.

In the following theorem, we show that after unlearning, our Targeted Sample Unlearning yields a more concentrated parameter update compared to the application of gradient ascent.

THEOREM 4.1 (PARAMETER UPDATES FOR δ -TARGETED ^{λ}). Consider a pre-trained neural network model with parameters $\hat{\theta} = \hat{\psi}, \hat{\phi}$, where $\hat{\psi}$ represents the parameter of all layers except the last fully connected layer, $\hat{\phi}$ represents the parameter of the last fully connected layer, C is the number of classes, and M is the dimension of the feature representation. Let $\mathbf{o} = g_{\hat{\psi}}(\mathbf{x}) \in \mathbb{R}^M$ represent the activation before the last linear layer for an input \mathbf{x} , with all components $o_i > 0$ (e.g., after ReLU activation). The logits are given by $f_{\hat{\theta}}(\mathbf{x}) = h_{\hat{\phi}}(\mathbf{o}) = \hat{\phi}\mathbf{o}$. We compare two updated parameter sets ϕ'' and ϕ' obtained after unlearning a to-be-forgotten example (\mathbf{x}, y) :

- 1. Our δ -Targeted ^{λ} Method Update (TD) with Pseudo-Label y' :

$$\phi'' = \hat{\phi} - \eta \nabla_{\hat{\phi}} \ell(\sigma(h_{\hat{\phi}}(\mathbf{o})), \mathbf{y}') = \hat{\phi} - \eta (\sigma(h_{\hat{\phi}}(\mathbf{o})) - \mathbf{y}') \mathbf{o}^\top;$$

- 2. Gradient Ascent Update (GA) with Original Label y :

$$\phi' = \hat{\phi} + \eta \nabla_{\hat{\phi}} \ell(\sigma(h_{\hat{\phi}}(\mathbf{o})), \mathbf{y}) = \hat{\phi} + \eta' (\sigma(h_{\hat{\phi}}(\mathbf{o})) - \mathbf{y}) \mathbf{o}^\top,$$

where σ is the softmax function, ℓ is the cross-entropy loss, \mathbf{y} is the one-hot encoding of the original label, \mathbf{y}' is the one-hot encoding of the pseudo-label assigned by our method, and η and η' are the learning rates. According to Definition 4.1, the parameter updates in ϕ'' are more concentrated compared to those in ϕ' . Specifically, the updates in ϕ'' are concentrated on the weights associated with the pseudo-label y' and the original label y , while the updates in ϕ' are distributed across all classes.

Proof Intuition. In gradient ascent, the update direction is the gradient of the loss with respect to the original label \mathbf{y} . This gradient changes all parameters associated with not only the true class y but also all other classes, leading to more uniformly distributed parameter changes. In our method, we perform gradient descent

²An intuitive explanation for this phenomenon is that changes to the last fully connected layer's parameters get amplified and propagated to former layers due to the chain rule employed during backpropagation.

on the loss with respect to the pseudo-label y' , which is the second most confident prediction. The gradient primarily affects the parameters associated with the pseudo-label y' and the activations corresponding to o . Since y' is chosen to be the class the model is already somewhat confident about, the parameter changes are concentrated on the weights connecting o to y' , leading to more focused parameter updates. Therefore, under the condition that the total magnitude of changes is equal, our method results in parameter changes that are more concentrated on a subset of parameters, which satisfies the concentrated parameter update in Definition 4.1. Proof details can be found in the appendix.

Theorem 4.1 provides a theoretical justification for the δ -Targeted¹ method. It is important to note that for our theoretical analysis, we set both $\lambda = 1$ and $\delta = 1$. The parameter δ controls the privacy level. Setting $\delta = 1$ ensures that the model mispredicts labels for all to-be-forgotten examples after fine-tuning. We chose $\delta = 1$ specifically to analyze how parameter changes induced by the δ -Targeted¹ method compare to the method induced by a gradient ascent when a model is intensively fine-tuned to force misprediction of an example. Additionally, we employ a first-order cross-entropy loss. By increasing the order of the cross-entropy loss (i.e., setting $\lambda > 1$), the change to the model should be more focused. This is because higher values of λ emphasize larger losses, causing the network to concentrate more on examples with large training losses. Empirically, results for $\lambda = 1$ and $\lambda = 3$ are presented in Section 5.1. The results are consistent where $\lambda = 3$ performs better in most cases when the size of to-be-forgotten examples is large.

5 Experiments

Datasets and Baselines The evaluation of our unlearning methods was conducted using four different benchmark image classification datasets: CIFAR-10, CIFAR-100 [18], Fashion-MNIST [34], and Tiny-ImageNet [20]. CIFAR-10 contains 10 classes with 50,000 training images and 10,000 test images. CIFAR-100 includes 100 classes with 50,000 training images and 10,000 test images. Fashion-MNIST has 10 classes with 60,000 training images and 10,000 test images. Tiny-ImageNet contains 200 classes with 100,000 training images and 10,000 test images.

We utilized ResNet-18, ResNet-50 [14] and vision transformer (ViT) architectures [7] for our experiments. We established a comparison of our proposed methods against different baselines. These include “BEFORE”, which represents the state of the model before unlearning; “NegGrad”, which involves fine-tuning models using gradient ascent [19]; “RandL”, which misleads the model through random label perturbation [15]; “Bad Teacher”, which uses competent and incompetent teachers in a student-teacher framework to induce forgetfulness [6]; “Amnesiac”, which undoes the parameter updates from only the batches containing sensitive data [12]. Our proposed methods are δ -Targeted¹, which sets the $\lambda = 1$ and δ -Targeted³, which sets $\lambda = 3$. Note that for “Bad Teacher” and “Amnesiac”, some remaining examples need to be accessed. In our setting, we provide them with 300 remaining examples. For all other methods, no remaining examples are provided.

Experiment Settings For each dataset, a random subset of images from the training dataset was selected to form the unlearning sample, denoted as D_f . The size of this sample varied, with cardinalities

set at 4, 16, 64, 128, 256, 512, 1024, respectively. For our method, the unlearning process was carried out using an Adam optimizer. The parameters of the Adam optimizer were kept constant across all experiments, with a learning rate of $1e - 4$, a weight decay of $1e - 5$. To ensure statistical significance, all experiments were conducted five times, each with different random initializations for each setting. The standard deviation of the results from each five trials has also been reported.

Note that to make the pre-trained model randomly guess a to-be-forgotten example, δ should be set to $1/C$, where C is the number of classes. In our experiments, we test our method under the most challenging conditions, we set $\delta = 0$ in our experiments. This requires the model to misclassify all to-be-forgotten examples by learning from the label noise while preserving performance on the remaining examples.

5.1 Magnitude of Parameter Changes for Different Methods

In Section 4, we showed that our method results in more concentrated parameter updates for the last-layer parameters than other methods. We hypothesize that concentrated updates in the last-layer parameters should lead to smaller changes in all parameters. We validate this on different neural network architectures, datasets, and baselines.

In Fig. 1, we measure and compare the parameter changes introduced by different unlearning methods, quantified using the L_1 -norm. The results confirm that our proposed method consistently induces smaller parameter changes after unlearning. This empirical evidence supports our method’s effectiveness in not only forgetting specific data but also in inducing only small changes to the model’s parameters. Moreover, the consistent pattern of small parameter changes observed across various datasets and neural network architectures underscores the robustness of our method across different datasets and models.

5.2 Relations between Accuracies and Parameter Changes

In Fig. 2, we visualize the relationship between the change in all parameters of the model and test accuracy using scatter plots. The change in the models’ parameters is quantified using the L_1 -norm. The dependence between changes in last-layer parameters and model accuracy on remaining examples using our method is illustrated in App.B.9.

For each subplot, δ -Targeted¹ and δ -Targeted³ are collectively displayed. A linear regression line obtained by regressing models’ test accuracies against changes in parameters is also shown. The results indicate that there is a consistent negative dependence between the change in the models’ parameters $\Delta\theta$ and test accuracy across various datasets. These findings suggest that small changes in parameters during the unlearning process contribute to maintaining test accuracy. This confirms that our method, by inducing small changes to the model’s parameters, successfully maintains the performance of the model after unlearning.

Additionally, smaller models seem to be more sensitive to parameter changes, which impacts their generalization ability. By examining Fig. 2 and Fig. 7 in the appendix, it is evident that across

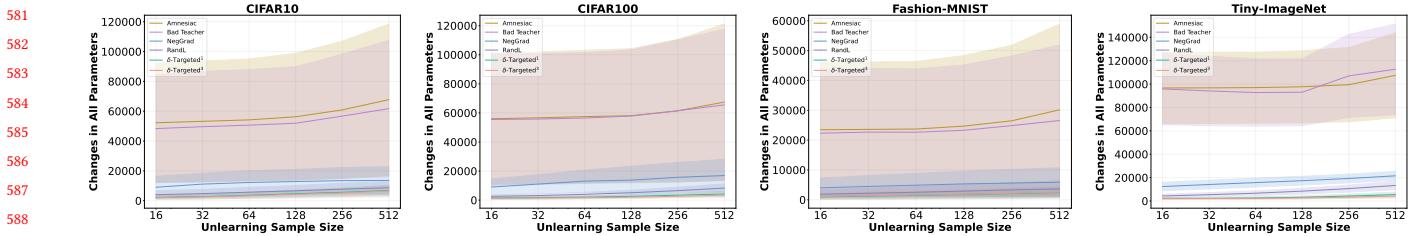


Figure 1: Parameter changes for different numbers of to-be-forgotten examples across different datasets. The figure includes three neural network architectures: ResNet-18, ResNet-50, and ViT. The changes are quantified by the L_1 -norm. Our method results in the smallest changes to the models' parameters across the datasets. Results for individual neural network architectures are provided in App.B.6.

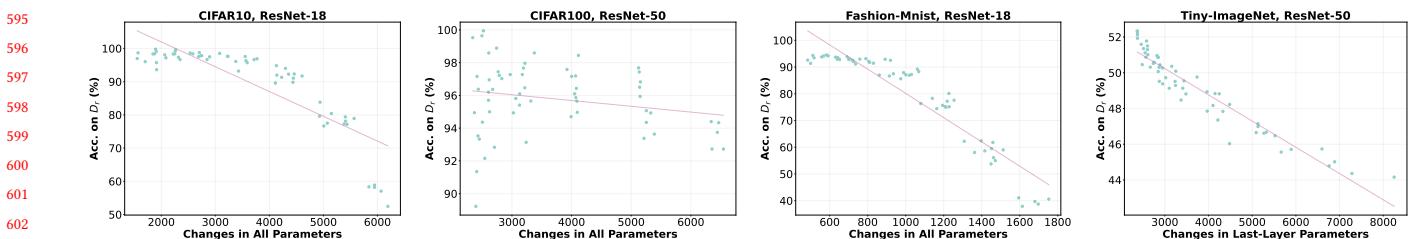


Figure 2: Dependence between changes in all parameters and model accuracy on remaining examples by using our method. The result shows a negative correlation. The results show a positive correlation. Results for other neural network architectures are provided in App.B.8.

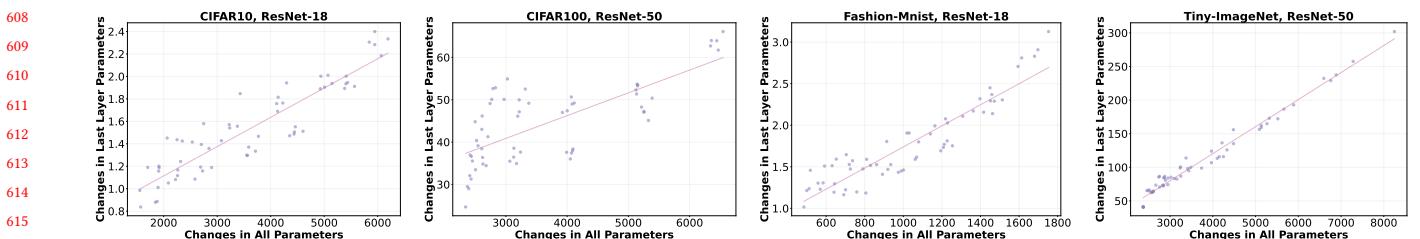


Figure 3: Dependence between the changes in all parameters and the changes in the last layer parameters by using our method. The magnitude of the parameter changes is quantified by the L_1 -norm. The results show a positive correlation. Results for other neural network architectures are provided in App.B.7.

most datasets, the same amount of parameter change causes a larger drop in accuracy for smaller neural network models. Another interesting observation from the figure is that the rate at which accuracy decreases in response to parameter changes varies by dataset. The accuracy drop is fastest on Fashion-MNIST, followed by CIFAR-10, CIFAR-100, and Tiny-ImageNet. This trend may suggest that models trained on more complex tasks are not as vulnerable to losses in accuracy due to parameter changes.

5.3 Relations between Changes in Last-Layer Parameters and Changes in All Parameters

In Section 4, we showed that our method results in concentrated parameter updates in the last-layer parameters. We focused on parameter changes in the last layer instead of changes in all parameters of a model. We hypothesize that updates in the last-layer

parameters should lead to smaller changes in all parameters. We validate this on different neural network architectures and datasets in this subsection.

Fig. 3 demonstrates the relationship between changes in the last-layer parameters and changes in a model's overall parameters. Each subplot represents a dataset and neural network architecture combination. The data reveal a consistent positive correlation between changes in the last-layer parameters and changes in a model's overall parameters. This validated our guess that changes in the last layer's parameters can effectively represent changes across all parameters of a model.

5.4 Performance on Different Datasets

Tables 1, 2, 3, and 4 present the classification accuracies on different datasets with different methods for various sizes of the unlearning

Table 1: Means and standard deviations (percentage) of classification accuracy on CIFAR10.

Method	$ D_f =16$	$ D_f =32$	$ D_f =64$	$ D_f =128$	$ D_f =256$	$ D_f =512$
D_f	<i>BEFORE</i>	97.92 \pm 4.93	98.12 \pm 3.39	98.33 \pm 3.15	98.39 \pm 2.73	98.26 \pm 2.79
	NegGrad	16.25 \pm 33.68	3.33 \pm 4.90	2.92 \pm 4.30	3.12 \pm 4.60	3.62 \pm 5.19
	RandL	0.00 \pm 0.00				
	Bad Teacher	11.25 \pm 8.29	9.38 \pm 4.70	11.46 \pm 3.68	12.34 \pm 4.64	13.18 \pm 4.34
	Amnesiac	0.00 \pm 0.00				
	δ -Targeted ¹	0.00 \pm 0.00				
D_r	δ -Targeted ³	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.18 \pm 0.35
	<i>BEFORE</i>	97.88 \pm 3.00	97.88 \pm 3.00	97.88 \pm 3.00	97.87 \pm 3.00	97.88 \pm 2.99
	NegGrad	69.62 \pm 22.88	53.10 \pm 24.18	32.81 \pm 15.36	17.46 \pm 6.98	9.29 \pm 1.36
	RandL	87.98 \pm 6.54	84.92 \pm 7.67	74.79 \pm 10.21	59.97 \pm 9.95	42.47 \pm 9.91
	Bad Teacher	78.09 \pm 6.95	78.11 \pm 7.25	75.75 \pm 7.26	69.23 \pm 10.53	56.55 \pm 14.55
	Amnesiac	76.83 \pm 4.44	75.84 \pm 4.27	73.77 \pm 3.97	68.89 \pm 4.39	54.21 \pm 7.68
D_t	δ -Targeted ¹	94.35 \pm 6.48	93.73 \pm 7.13	89.34 \pm 10.51	84.57 \pm 10.22	75.32 \pm 7.38
	δ -Targeted ³	93.24 \pm 4.51	93.61 \pm 5.50	91.98 \pm 7.80	90.55 \pm 8.50	86.22 \pm 7.98
	<i>BEFORE</i>	89.62 \pm 5.62				
	NegGrad	64.41 \pm 20.62	48.77 \pm 21.49	29.71 \pm 13.12	15.61 \pm 5.62	8.70 \pm 1.45
	RandL	80.15 \pm 7.11	76.86 \pm 7.97	67.25 \pm 9.89	53.41 \pm 8.61	37.44 \pm 8.49
	Bad Teacher	73.53 \pm 4.35	73.29 \pm 5.23	71.26 \pm 5.47	65.24 \pm 9.19	53.46 \pm 13.09

Table 2: Means and standard deviations (percentage) of classification accuracy on CIFAR100.

Method	$ D_f =16$	$ D_f =32$	$ D_f =64$	$ D_f =128$	$ D_f =256$	$ D_f =512$
D_f	<i>BEFORE</i>	95.42 \pm 8.06	93.96 \pm 9.98	94.48 \pm 7.94	94.84 \pm 7.19	94.71 \pm 7.38
	NegGrad	0.42 \pm 1.56	0.21 \pm 0.78	0.21 \pm 0.53	7.03 \pm 24.85	0.26 \pm 0.47
	RandL	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00
	Bad Teacher	1.25 \pm 3.39	1.25 \pm 1.91	2.08 \pm 2.03	2.24 \pm 2.23	3.12 \pm 2.58
	Amnesiac	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00
	δ -Targeted ¹	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00
D_r	δ -Targeted ³	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00
	<i>BEFORE</i>	95.04 \pm 6.97	95.04 \pm 6.97	95.04 \pm 6.96	95.04 \pm 6.97	95.04 \pm 6.97
	NegGrad	62.90 \pm 19.23	59.81 \pm 25.59	53.14 \pm 29.77	47.59 \pm 34.69	41.18 \pm 29.50
	RandL	82.07 \pm 11.83	83.30 \pm 14.32	81.64 \pm 17.73	78.32 \pm 22.15	73.81 \pm 26.68
	Bad Teacher	44.56 \pm 21.53	43.94 \pm 20.06	40.29 \pm 17.11	31.14 \pm 17.08	23.36 \pm 14.10
	Amnesiac	41.52 \pm 16.65	40.75 \pm 15.86	38.51 \pm 14.30	35.55 \pm 13.66	26.92 \pm 14.24
D_t	δ -Targeted ¹	93.39 \pm 6.96	92.85 \pm 7.43	92.00 \pm 8.63	90.64 \pm 10.10	89.32 \pm 11.06
	δ -Targeted ³	90.88 \pm 5.42	91.58 \pm 6.31	91.09 \pm 7.33	90.28 \pm 8.20	89.70 \pm 8.62
	<i>BEFORE</i>	67.07 \pm 8.89	67.07 \pm 8.89	67.07 \pm 8.89	67.07 \pm 8.89	67.07 \pm 8.89
	NegGrad	43.66 \pm 11.78	41.15 \pm 15.74	36.29 \pm 18.85	32.19 \pm 23.34	27.50 \pm 19.32
	RandL	55.43 \pm 7.83	56.48 \pm 9.79	55.72 \pm 12.13	53.40 \pm 14.77	49.57 \pm 17.30
	Bad Teacher	34.26 \pm 12.22	33.89 \pm 11.57	31.59 \pm 10.02	24.47 \pm 11.97	18.69 \pm 10.84

sample D_f . Note that the results in these tables are obtained by averaging over different neural network architectures, including ResNet-18, ResNet-50, and ViT. The results for individual neural network architectures can be found in App. B.

We first examine the performance on D_f , which is the subset of data we want our model to unlearn. The tables show that all methods can effectively unlearn the data, resulting in near-zero accuracy on D_f . Secondly, looking at the performance on D_r , which is the remaining examples that the model should still remember, it is evident that both δ -Targeted¹ and δ -Targeted³ consistently outperform the other methods. In particular, δ -Targeted³ demonstrates superior

performance compared to other methods, as it mostly achieves the highest accuracy across different sizes of D_f , indicating that it effectively retains the model's ability to predict the remaining sample despite the unlearning process.

Lastly, on D_t , the test dataset, both δ -Targeted¹ and δ -Targeted³ again demonstrate superior performance. Despite the unlearning process, these methods ensure the model maintains its overall predictive performance on unseen data, outperforming other methods across different sizes of D_f .

When the size of the unlearning sample D_f is large, δ -Targeted³ consistently outperforms the other methods. It effectively unlearns

Table 3: Means and standard deviations (percentage) of classification accuracy on Fashion-MNIST.

	Method	$ D_f =16$	$ D_f =32$	$ D_f =64$	$ D_f =128$	$ D_f =256$	$ D_f =512$
D_f	BEFORE	96.25 \pm 5.00	95.42 \pm 4.25	94.79 \pm 3.16	94.95 \pm 2.73	95.13 \pm 2.11	95.76 \pm 1.78
	NegGrad	2.92 \pm 6.80	2.71 \pm 6.54	3.65 \pm 6.42	2.24 \pm 4.33	3.36 \pm 5.03	3.48 \pm 4.99
	RandL	0.00 \pm 0.00	0.46 \pm 0.68				
	Bad Teacher	7.92 \pm 8.06	6.04 \pm 5.41	6.77 \pm 5.84	6.93 \pm 5.51	8.83 \pm 6.67	8.87 \pm 7.05
	Amnesiac	0.00 \pm 0.00	0.42 \pm 1.56	2.50 \pm 3.41	6.30 \pm 8.82	12.37 \pm 13.05	10.68 \pm 10.84
	δ -Targeted ¹	0.00 \pm 0.00	0.00 \pm 0.00	0.62 \pm 1.11	1.09 \pm 1.91	4.56 \pm 6.83	17.50 \pm 24.87
	δ -Targeted ³	11.25 \pm 16.65	10.42 \pm 15.21	15.62 \pm 22.37	22.40 \pm 31.67	25.55 \pm 36.15	25.85 \pm 36.57
D_r	BEFORE	96.04 \pm 1.37	96.04 \pm 1.37	96.04 \pm 1.37	96.04 \pm 1.36	96.04 \pm 1.36	96.04 \pm 1.36
	NegGrad	53.98 \pm 24.96	41.26 \pm 22.91	30.46 \pm 14.70	17.86 \pm 7.91	11.69 \pm 3.20	8.02 \pm 2.19
	RandL	72.05 \pm 17.74	64.08 \pm 23.78	53.59 \pm 24.06	40.22 \pm 19.96	27.35 \pm 14.38	18.87 \pm 8.96
	Bad Teacher	79.16 \pm 10.05	76.20 \pm 12.20	71.15 \pm 15.48	63.17 \pm 17.87	49.02 \pm 17.18	37.71 \pm 14.54
	Amnesiac	81.51 \pm 6.08	80.15 \pm 5.16	74.83 \pm 7.61	69.01 \pm 8.52	56.34 \pm 6.37	40.00 \pm 3.08
	δ -Targeted ¹	91.76 \pm 3.26	88.56 \pm 3.31	82.73 \pm 3.38	72.52 \pm 3.96	57.43 \pm 5.38	45.50 \pm 11.91
	δ -Targeted ³	92.47 \pm 1.68	91.98 \pm 2.15	90.04 \pm 2.48	85.04 \pm 2.44	75.70 \pm 4.92	61.69 \pm 11.73
D_t	BEFORE	91.95 \pm 0.33					
	NegGrad	51.19 \pm 23.07	38.70 \pm 21.18	28.15 \pm 13.11	16.12 \pm 6.77	10.58 \pm 2.58	7.39 \pm 2.37
	RandL	68.66 \pm 16.13	60.82 \pm 21.88	50.34 \pm 22.10	37.26 \pm 18.06	24.65 \pm 12.58	16.47 \pm 7.52
	Bad Teacher	77.94 \pm 10.06	74.83 \pm 12.05	69.89 \pm 15.42	61.94 \pm 17.84	47.99 \pm 17.12	37.14 \pm 14.48
	Amnesiac	80.24 \pm 6.23	78.87 \pm 5.31	73.62 \pm 7.83	67.77 \pm 8.67	55.18 \pm 6.86	38.94 \pm 3.54
	δ -Targeted ¹	87.99 \pm 1.62	84.85 \pm 2.47	79.04 \pm 3.16	68.78 \pm 4.71	54.04 \pm 6.72	42.59 \pm 13.57
	δ -Targeted ³	88.89 \pm 0.85	88.34 \pm 0.75	86.27 \pm 1.27	81.36 \pm 2.64	72.14 \pm 6.60	58.54 \pm 13.48

Table 4: Means and standard deviations (percentage) of classification accuracy on Tiny-ImageNet.

	Method	$ D_f =16$	$ D_f =32$	$ D_f =64$	$ D_f =128$	$ D_f =256$	$ D_f =512$
D_f	BEFORE	64.53 \pm 24.87	65.17 \pm 24.07	64.59 \pm 23.62	65.13 \pm 23.10	65.32 \pm 22.02	66.52 \pm 23.84
	NegGrad	0.00 \pm 0.00					
	RandL	0.00 \pm 0.00					
	Bad Teacher	0.42 \pm 1.56	0.21 \pm 0.78	0.00 \pm 0.00	0.42 \pm 0.48	0.31 \pm 0.36	0.42 \pm 0.40
	Amnesiac	0.00 \pm 0.00					
	δ -Targeted ¹	0.00 \pm 0.00					
	δ -Targeted ³	0.00 \pm 0.00	0.01 \pm 0.05				
D_r	BEFORE	71.22 \pm 15.44	73.26 \pm 15.19				
	NegGrad	44.78 \pm 9.53	48.91 \pm 9.60	49.47 \pm 10.10	43.80 \pm 11.20	38.23 \pm 11.02	34.23 \pm 8.57
	RandL	46.29 \pm 11.93	47.89 \pm 15.37	50.41 \pm 17.50	52.53 \pm 17.65	54.07 \pm 15.93	53.48 \pm 14.17
	Bad Teacher	25.22 \pm 4.70	25.55 \pm 5.01	24.58 \pm 3.91	23.25 \pm 5.12	19.22 \pm 5.42	16.43 \pm 6.06
	Amnesiac	23.86 \pm 8.20	23.34 \pm 7.86	23.03 \pm 7.56	22.90 \pm 8.15	22.51 \pm 7.88	19.23 \pm 8.16
	δ -Targeted ¹	61.68 \pm 15.20	61.02 \pm 16.13	59.60 \pm 17.15	59.19 \pm 17.17	59.62 \pm 16.62	57.34 \pm 17.32
	δ -Targeted ³	60.84 \pm 12.44	59.67 \pm 13.10	60.13 \pm 14.65	60.05 \pm 15.72	60.04 \pm 16.49	59.98 \pm 16.08
D_t	BEFORE	52.51 \pm 0.23	52.51 \pm 0.23	52.51 \pm 0.23	52.51 \pm 0.23	52.53 \pm 0.24	52.53 \pm 0.24
	NegGrad	39.32 \pm 2.95	40.42 \pm 2.47	40.13 \pm 2.42	38.30 \pm 2.62	32.08 \pm 3.72	30.09 \pm 2.27
	RandL	40.37 \pm 3.05	40.58 \pm 4.67	40.88 \pm 5.49	42.48 \pm 4.41	43.67 \pm 3.74	44.17 \pm 3.28
	Bad Teacher	24.66 \pm 5.75	25.35 \pm 5.11	26.43 \pm 6.38	22.85 \pm 5.12	18.95 \pm 5.66	16.92 \pm 6.53
	Amnesiac	22.98 \pm 7.84	24.58 \pm 9.00	23.35 \pm 8.36	24.21 \pm 9.51	23.13 \pm 9.24	17.96 \pm 8.23
	δ -Targeted ¹	51.28 \pm 1.30	50.68 \pm 2.18	49.79 \pm 4.38	48.59 \pm 2.98	47.97 \pm 3.28	46.90 \pm 3.68
	δ -Targeted ³	50.04 \pm 2.00	50.39 \pm 1.17	49.84 \pm 3.36	50.77 \pm 2.69	49.37 \pm 2.29	47.02 \pm 3.74

harder-to-forget examples by emphasizing the importance of harder-to-forget examples in the unlearning set. This makes it particularly advantageous for larger unlearning sets. However, when the unlearning set is small, δ -Targeted¹ tends to outperform δ -Targeted³. This might be because when all examples can be easily unlearned, placing extra emphasis on certain examples (as in $\lambda = 3$) does not provide additional benefits. A potential reason could be that, when the to-be-forgotten sample size is small, these examples can be unlearned in just a few steps. Emphasizing the loss too much at the beginning can lead to large parameter updates, which adversely affect the model's performance on the remaining data.

6 Conclusion

In this work, we presented δ -Targeted^Y, a method designed to retain model performance after removing a specific training sample without the need for reassessing the remaining sample. Theoretically, we demonstrated that our method results in more concentrated parameter updates than those introduced by gradient ascent, thereby effectively reducing the performance drops commonly observed after unlearning. Our empirical evaluations further corroborate these findings, as δ -Targeted^Y consistently achieves state-of-the-art results across a variety of benchmark image datasets, underscoring both its theoretical soundness and practical efficacy.

References

- [1] Jonathan Brophy and Daniel Lowd. 2021. Machine unlearning for random forests. In *International Conference on Machine Learning*. PMLR, 1092–1104.
- [2] Yinzhi Cao and Junfeng Yang. 2015. Towards making systems forget with machine unlearning. In *2015 IEEE Symposium on Security and Privacy*. IEEE, 463–480.
- [3] Abraham Chan, Arpan Gujarati, Karthik Pattabiraman, and Sathish Gopalakrishnan. [n. d.]. Hierarchical Unlearning Framework for Multi-Class Classification. In *NeurIPS 2024 Workshop on Fine-Tuning in Modern Machine Learning: Principles and Scalability*.
- [4] Min Chen, Weizhuo Gao, Gaoyang Liu, Kai Peng, and Chen Wang. 2023. Boundary Unlearning: Rapid Forgetting of Deep Networks via Shifting the Decision Boundary. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 7766–7775.
- [5] Vikram S Chundawat, Ayush K Tarun, Murari Mandal, and Mohan Kankanhalli. 2022. Zero-shot machine unlearning. *arXiv preprint arXiv:2201.05629* (2022).
- [6] Vikram S Chundawat, Ayush K Tarun, Murari Mandal, and Mohan Kankanhalli. 2023. Can bad teaching induce forgetting? unlearning in deep networks using an incompetent teacher. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 37. 7210–7217.
- [7] Alexey Dosovitskiy. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020).
- [8] Sayna Ebrahimi, Franziska Meier, Roberto Calandra, Trevor Darrell, and Marcus Rohrbach. 2020. Adversarial continual learning. In *European Conference on Computer Vision*. Springer, 386–402.
- [9] Jack Foster, Stefan Schoepf, and Alexandra Brintrup. 2024. Fast machine unlearning without retraining through selective synaptic dampening. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 12043–12051.
- [10] Antonio Ginart, Melody Guan, Gregory Valiant, and James Y Zou. 2019. Making ai forget you: Data deletion in machine learning. *Advances in neural information processing systems* 32 (2019).
- [11] Aditya Golatkar, Alessandro Achille, and Stefano Soatto. 2020. Eternal sunshine of the spotless net: Selective forgetting in deep networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9304–9312.
- [12] Laura Graves, Vineel Nagisetty, and Vijay Ganesh. 2021. Amnesia machine learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 11516–11524.
- [13] Ling Han, Nanqing Luo, Hao Huang, Jing Chen, and Mary-Anne Hartley. 2024. Towards Independence Criterion in Machine Unlearning of Features and Labels. *arXiv preprint arXiv:2403.08124* (2024).
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [15] Mark He Huang, Lin Geng Foo, and Jun Liu. 2024. Learning to unlearn for robust machine unlearning. *arXiv preprint arXiv:2407.10494* (2024).
- [16] Jinghan Jia, Jiancheng Liu, Parikshit Ram, Yuguang Yao, Gaowen Liu, Yang Liu, Pranay Sharma, and Sijia Liu. 2023. Model sparsification can simplify machine unlearning. *arXiv preprint arXiv:2304.04934* (2023).
- [17] Junyaup Kim and Simon S Woo. 2022. Efficient Two-Stage Model Retraining for Machine Unlearning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4361–4369.
- [18] Alex Krizhevsky, Geoffrey Hinton, et al. 2009. Learning multiple layers of features from tiny images. (2009).
- [19] Meghdad Kurmanji, Peter Triantafillou, Jamie Hayes, and Eleni Triantafillou. 2024. Towards unbounded machine unlearning. *Advances in neural information processing systems* 36 (2024).
- [20] Ya Le and Xuan Yang. 2015. Tiny imagenet visual recognition challenge. *CS 231N* 7, 7 (2015), 3.
- [21] Sangyong Lee and Simon S Woo. 2023. Undo: Effective and accurate unlearning method for deep neural networks. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*. 4043–4047.
- [22] Guanghao Li, Li Shen, Yan Sun, Yue Hu, Han Hu, and Dacheng Tao. 2023. Subspace based Federated Unlearning. *arXiv preprint arXiv:2302.12448* (2023).
- [23] Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, and Xingjun Ma. 2021. Anti-backdoor learning: Training clean models on poisoned data. *Advances in Neural Information Processing Systems* 34 (2021), 14900–14912.
- [24] Ananth Mahadevan and Michael Mathioudakis. 2021. Certifiable machine unlearning for linear models. *arXiv preprint arXiv:2106.15093* (2021).
- [25] Ronak Mehta, Sourav Pal, Vikas Singh, and Sathy N Ravi. 2022. Deep Unlearning via Randomized Conditionally Independent Hessians. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10422–10431.
- [26] Samuele Poppi, Sara Sarto, Marcella Cornia, Lorenzo Baraldi, and Rita Cucchiara. 2024. Multi-Class Unlearning for Image Classification via Weight Filtering. *IEEE Intelligent Systems* (2024).
- [27] General Data Protection Regulation. 2018. General data protection regulation (GDPR). *Intersoft Consulting*, Accessed in October 24, 1 (2018).
- [28] Amrit Sethur, Benjamin Eysenbach, Virginia Smith, and Sergey Levine. 2022. Adversarial unlearning: Reducing confidence along adversarial directions. *Advances in Neural Information Processing Systems* 35 (2022), 18556–18570.
- [29] Takashi Shibata, Go Irie, Daiki Ikami, and Yu Mitsuzumi. 2021. Learning with Selective Forgetting.. In *IJCAI*, Vol. 3. 4.
- [30] Vinith M. Suriyakumar and Ashia C. Wilson. 2022. Algorithms that Approximate Data Removal: New Results and Limitations.
- [31] Ayush K Tarun, Vikram S Chundawat, Murari Mandal, and Mohan Kankanhalli. 2021. Fast yet effective machine unlearning. *arXiv preprint arXiv:2111.08947* (2021).
- [32] Ayush K Tarun, Vikram S Chundawat, Murari Mandal, and Mohan Kankanhalli. 2023. Fast yet effective machine unlearning. *IEEE Transactions on Neural Networks and Learning Systems* (2023).
- [33] Leijie Wu, Song Guo, Junxiao Wang, Zicong Hong, Jie Zhang, and Yaohong Ding. 2022. Federated unlearning: Guarantee the right of clients to forget. *IEEE Network* 36, 5 (2022), 129–135.
- [34] Han Xiao, Kashif Rasul, and Roland Vollgraf. 2017. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747* (2017).
- [35] Jingwen Ye, Yifang Fu, Jie Song, Xingyi Yang, Songhua Liu, Xin Jin, Mingli Song, and Xinchao Wang. 2022. Learning with Recoverable Forgetting. *arXiv preprint arXiv:2207.08224* (2022).
- [36] Youngsik Yoon, Jinhyeon Nam, Hyojeong Yun, Dongwoo Kim, and Jungseul Ok. 2022. Few-Shot Unlearning by Model Inversion. *arXiv preprint arXiv:2205.15567* (2022).
- [37] Siyao Yu, Fei Sun, Jiafeng Guo, Ruqing Zhang, and Xueqi Cheng. 2022. Legonet: A fast and exact unlearning architecture. *arXiv preprint arXiv:2210.16023* (2022).

929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986

987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044

1045 A Proof of Theorem 4.1

1046 We define some notations used in the proof first. Let (\mathbf{x}, y) by an example. Let \mathbf{o} represent the activation just before the last layer of the
 1047 pre-trained model. For the label y , its corresponding one-hot encoding is denoted by \mathbf{y} . The estimated class posteriors $\mathbf{p} = \sigma \circ h_{\hat{\phi}}(\mathbf{o})$ are
 1048 defined by the output of the softmax function applied to the linear transformation. It can be further expanded to $\mathbf{p} = \sigma(\hat{\phi}\mathbf{o}) = \sigma(\mathbf{z})$, where
 1049 $\hat{\phi} \in \mathbb{R}^{C \times M}$ represents the parameter of the last fully connected layer, $\mathbf{z} = \hat{\phi}\mathbf{o}$ represents the input to the softmax function, and σ is the
 1050 softmax function.
 1051

1052 Firstly, we derive the gradient of cross-entropy loss with respect to the last layer parameter ϕ . The cross-entropy loss function $\ell(\mathbf{p}, \mathbf{y})$ is
 1053 defined as:
 1054

$$\ell(\mathbf{p}, \mathbf{y}) = - \sum_{k=1}^C \mathbf{y}_k \cdot \log(\mathbf{p}_k). \quad (3)$$

1056 The gradient of ℓ with respect to the predicted class-posterior probability $\mathbf{p}_k = P(Y = k | X = \mathbf{x})$ can be calculated as:
 1057

$$\frac{\partial \ell}{\partial \mathbf{p}_k} = - \frac{\mathbf{y}_k}{\mathbf{p}_k}. \quad (4)$$

1060 Consider applying softmax function $\sigma(z)_i$ to calculate i -th class posterior $\mathbf{p}_i = \sigma(z)_i = \frac{e^{z_i}}{\sum_{j=1}^C e^{z_j}}$.
 1061

1062 Case when $i = k$, the derivative $\frac{\partial \mathbf{p}_i}{\partial z_k} = \frac{\partial \mathbf{p}_k}{\partial z_k}$ is
 1063

$$\begin{aligned} \frac{\partial \mathbf{p}_k}{\partial z_k} &= \frac{\sum_{j=1}^C e^{z_j} \cdot e^{z_k} - e^{z_k} \cdot e^{z_k}}{(\sum_{j=1}^C e^{z_j})^2} \\ &= \frac{e^{z_k}}{\sum_{j=1}^C e^{z_j}} \cdot \left(1 - \frac{e^{z_k}}{\sum_{j=1}^C e^{z_j}}\right) \\ &= \mathbf{p}_k(1 - \mathbf{p}_k). \end{aligned}$$

1071 Case when $i \neq k$, the derivative $\frac{\partial \mathbf{p}_i}{\partial z_k}$ is
 1072

$$\begin{aligned} \frac{\partial \mathbf{p}_i}{\partial z_k} &= \frac{0 - e^{z_k} \cdot e^{z_i}}{(\sum_{j=1}^C e^{z_j})^2} \\ &= - \frac{e^{z_k} \cdot e^{z_i}}{(\sum_{j=1}^C e^{z_j})^2} \\ &= - \mathbf{p}_k \mathbf{p}_i. \end{aligned}$$

1078 Applying the chain rule, the gradient of ℓ with respect to z_k becomes:
 1079

$$\frac{\partial \ell}{\partial z_k} = \sum_i \frac{\partial \ell}{\partial \mathbf{p}_i} \frac{\partial \mathbf{p}_i}{\partial z_k} = -\mathbf{y}_k(1 - \mathbf{p}_k) + \sum_{j \neq k} \mathbf{y}_j \mathbf{p}_j. \quad (5)$$

1084 Note that \mathbf{y} is the one-hot encoding form of y , \mathbf{y}_k can only include a value of either 0 or 1. The above equation simplifies to $\mathbf{p}_k - 1$ when
 1085 $y = k$ and \mathbf{p}_k when $y \neq k$. In matrix form, we can rewrite it as:
 1086

$$\frac{\partial \ell}{\partial \mathbf{z}} = \mathbf{p} - \mathbf{y}. \quad (6)$$

1089 Furthermore, the derivative of logits \mathbf{z} with respect to the last-layer parameter ϕ is $\frac{\partial \mathbf{z}}{\partial \phi} = \mathbf{o}^T$. Then, by applying the chain rule again, the
 1090 gradient of the loss ℓ with respect to the parameter ϕ is given by:
 1091

$$\frac{\partial \ell}{\partial \phi} = \frac{\partial \ell}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \phi} = (\mathbf{p} - \mathbf{y}) \mathbf{o}^T. \quad (7)$$

1095 To unlearn the example (\mathbf{x}, y) , we update ϕ to ϕ' using gradient ascent: $\phi' = \phi + \eta' \nabla_{\hat{\phi}} \ell(h_{\hat{\phi}}(\mathbf{o}), \mathbf{y})$, where $\eta' \in R$ is the learning rate, and
 1096 $\eta' > 0$. The updated elements of ϕ' are defined as:
 1097

$$\begin{aligned} \phi'_y &= \phi_y + \eta'(\mathbf{p}_y - 1)\mathbf{o}^T, \\ \phi'_k &= \phi_k + \eta'(\mathbf{p}_k)\mathbf{o}^T, \text{ for } k \neq y. \end{aligned}$$

1161 let $\Delta \mathbf{p} = \mathbf{p} - \mathbf{y}$, we write the above equations by using gradient ascent in matrix form

$$\begin{aligned}\phi' &= \phi + \eta'(\mathbf{p} - \mathbf{y})\mathbf{o}^T \\ &= \phi + \eta'\Delta\mathbf{p}\mathbf{o}^T.\end{aligned}$$

1165 We also consider an update of ϕ to ϕ'' using gradient descent: $\phi'' = \phi - \eta\nabla_{\hat{\phi}}L(h_{\hat{\phi}}(\mathbf{o}), y')$, where $\eta \in R$ is the learning rate, and $\eta > 0$.
 1166 The updated elements of ϕ'' are defined as:

$$\begin{aligned}\phi''_{y'} &= \phi_{y'} + \eta(1 - \mathbf{p}_{y'})\mathbf{o}^T, \\ \phi''_k &= \phi_k + \eta\mathbf{p}_k\mathbf{o}^T, \\ \phi''_y &= \phi_y + \eta\mathbf{p}_y\mathbf{o}^T, \text{ for } k \neq y, y'.\end{aligned}$$

1173 Let $\Delta\mathbf{p}' = \mathbf{y}' - \mathbf{p}$, we write the above equations by using gradient descent in matrix form

$$\begin{aligned}\phi'' &= \phi + \eta(\mathbf{y}' - \mathbf{p})\mathbf{o}^T \\ &= \phi + \eta\Delta\mathbf{p}'\mathbf{o}^T.\end{aligned}$$

1178 The L_1 -norms of $\Delta\phi'$ and $\Delta\phi''$ are given by:

$$\begin{aligned}\|\Delta\phi'\|_1 &= \eta' \|\Delta\mathbf{p}\mathbf{o}^T\|_1, \\ \|\Delta\phi''\|_1 &= \eta \|\Delta\mathbf{p}'\mathbf{o}^T\|_1.\end{aligned}$$

1182 We now compare the effect of gradient ascent and the gradient descent when $\Delta\phi'$ and $\Delta\phi''$ have the same L_1 norm. Recall $\Delta\mathbf{p} = \mathbf{p} - \mathbf{y}$
 1183 and $\Delta\mathbf{p}' = \mathbf{y}' - \mathbf{p}$, we can derive the following points:

- 1184 • For $i \neq y$ and $i \neq y'$, we have $|\Delta\mathbf{p}'_i| = |\mathbf{p}_i|$ and $|\Delta\mathbf{p}_i| = |\mathbf{p}_i|$.
- 1185 • When $i = y$, we have $|\Delta\mathbf{p}'_i| = |\mathbf{p}_y|$ and $|\Delta\mathbf{p}_i| = |\mathbf{p}_y - 1|$.
- 1186 • When $i = y'$, we have $|\Delta\mathbf{p}'_i| = |1 - \mathbf{p}_{y'}|$ and $|\Delta\mathbf{p}_i| = |-\mathbf{p}_{y'}|$.

1188 Therefore, the L_1 -norms of $\Delta\mathbf{p}$ and $\Delta\mathbf{p}'$ are:

$$\begin{aligned}\|\Delta\mathbf{p}\|_1 &= |\mathbf{p}_{y'}| + |1 - \mathbf{p}_y| + \sum_{i \neq y, y'} |\mathbf{p}_i| = 2 - 2\mathbf{p}_y, \\ \|\Delta\mathbf{p}'\|_1 &= |1 - \mathbf{p}_{y'}| + |\mathbf{p}_y| + \sum_{i \neq y, y'} |-\mathbf{p}_i| = 2 - 2\mathbf{p}_{y'}.\end{aligned}$$

1194 Given the fact that y is most confident label, then $\mathbf{p}_y > \mathbf{p}_{y'}$, we then have:

$$\begin{aligned}\|\Delta\mathbf{p}\|_1 - \|\Delta\mathbf{p}'\|_1 &= 2 - 2\mathbf{p}_y - (2 - 2\mathbf{p}_{y'}) \\ &= -2\mathbf{p}_y + 2\mathbf{p}_{y'} < 0.\end{aligned}$$

1198 Let $\eta' = \frac{1}{2-2\mathbf{p}_y}$ and $\eta = \frac{1}{2-2\mathbf{p}_{y'}}$ to ensure $\Delta\phi'$ and $\Delta\phi''$ have the same L_1 norm. Then, we can rewrite the update rules for ϕ' and ϕ'' as
 1199 follows:

$$\begin{aligned}\phi'_{*,j} &= \phi_{*,j} + \eta'(\mathbf{p} - \mathbf{y})\mathbf{o}_j = \phi_{*,j} + \frac{1}{2-2\mathbf{p}_y}(\mathbf{p} - \mathbf{y})\mathbf{o}_j, \\ \phi''_{*,j} &= \phi_{*,j} + \eta(\mathbf{y}' - \mathbf{p})\mathbf{o}_j = \phi_{*,j} + \frac{1}{2-2\mathbf{p}_{y'}}(\mathbf{y}' - \mathbf{p})\mathbf{o}_j,\end{aligned}$$

1206 where $\mathbf{M}_{*,j}$ denote the j -th column in a matrix \mathbf{M} .

1207 Now let us discuss the following two cases to compare the change of parameters:

1208 1). For all $i \neq y, y'$, we have:

$$\phi'_{ij} = \phi_{ij} + \frac{1}{2-2\mathbf{p}_y} \mathbf{p}_i \mathbf{o}_j, \quad (8)$$

$$\phi''_{ij} = \phi_{ij} - \frac{1}{2-2\mathbf{p}_{y'}} \mathbf{p}_i \mathbf{o}_j. \quad (9)$$

1215 Given that $\mathbf{p}_{y'} < \mathbf{p}_y$, it follows that $|\phi''_{ij}| - |\phi'_{ij}| < 0$. Thus, when the L_1 norm of $\Delta\phi'$ and $\Delta\phi''$ are the same, the gradient ascent method
 1216 leads to more changes in the parameters for predicting \mathbf{p}_i for $i \neq y$ and $i \neq y'$ compared to gradient descent.

1217 2). For $i = y$, we have:

$$\phi'_{ij} = \phi_{ij} + \frac{1}{2 - 2\mathbf{p}_y} (\mathbf{p}_y - 1) \mathbf{o}_j, \quad (10)$$

$$\phi''_{ij} = \phi_{ij} - \frac{1}{2 - 2\mathbf{p}_{y'}} \mathbf{p}_y \mathbf{o}_j. \quad (11)$$

Observing that $|\phi''_{ij}| - |\phi'_{ij}| = \frac{-1+\mathbf{p}_{y'}+\mathbf{p}_y}{2(1-\mathbf{p}_{y'})} < 0$ when having more than two classes. Then, when the L_1 norms of $\Delta\phi'$ and $\Delta\phi''$ are identical, the gradient ascent method leads to more changes in the parameters for predicting \mathbf{p}_y than gradient descent.

Combining the findings for both cases, $|\phi''_{ij}| - |\phi'_{ij}| < 0$. As L_1 norms of $\Delta\phi'$ and $\Delta\phi''$ are identical, we conclude that only for $i = y'$, the value of $|\phi''_{ij}| - |\phi'_{ij}| > 0$. This means our method only leads to large changes in the parameters for predicting pseudo-label y' .

When more than two classes, our method results in parameter updates that are more concentrated on the weights associated with the pseudo-label y' , while gradient ascent distributes updates more uniformly across all classes, which completes the proof.

B Additional Experiment Results

In this section, we provide additional experiment results, including performance on different datasets with various neural network structures, changes in last-layer parameters for different sample unlearning methods, and changes in all parameters for different sample unlearning methods. We also explore the relationships between changes in last-layer parameters and changes in all parameters, as well as the relationships between a model's accuracy on the remaining examples and changes in all parameters.

B.1 Performance on CIFAR10 with Different Neural Network Structures

Table 5: Means and standard deviations (percentage) of classification accuracy on CIFAR10 with ResNet-18.

	Method	$ D_f =16$	$ D_f =32$	$ D_f =64$	$ D_f =128$	$ D_f =256$	$ D_f =512$
D_f	BEFORE	100.00 \pm 0.00					
	NegGrad	40.00 \pm 48.99	0.00 \pm 0.00				
	RandL	0.00 \pm 0.00					
	Bad Teacher	7.50 \pm 7.29	10.62 \pm 4.24	11.56 \pm 2.90	9.53 \pm 2.39	10.23 \pm 1.72	11.48 \pm 1.18
	Amnesiac	0.00 \pm 0.00					
	δ -Targeted ¹	0.00 \pm 0.00					
	δ -Targeted ³	0.00 \pm 0.00					
D_r	BEFORE	99.99 \pm 0.00					
	NegGrad	92.52 \pm 7.43	74.05 \pm 10.62	46.34 \pm 7.70	25.25 \pm 3.48	9.28 \pm 1.58	4.97 \pm 0.84
	RandL	92.17 \pm 3.04	90.96 \pm 1.34	83.47 \pm 4.08	67.99 \pm 3.01	50.06 \pm 2.85	34.35 \pm 1.25
	Bad Teacher	77.73 \pm 4.66	81.09 \pm 1.33	81.18 \pm 1.27	77.40 \pm 1.62	67.22 \pm 4.14	45.04 \pm 6.04
	Amnesiac	78.64 \pm 1.34	79.30 \pm 1.26	78.10 \pm 1.46	74.23 \pm 1.23	60.68 \pm 1.34	35.88 \pm 2.96
	δ -Targeted ¹	98.39 \pm 1.05	98.83 \pm 0.59	96.61 \pm 1.88	92.33 \pm 1.88	79.59 \pm 2.51	56.98 \pm 2.33
	δ -Targeted ³	96.19 \pm 1.65	97.91 \pm 0.53	97.57 \pm 0.55	96.63 \pm 0.62	91.43 \pm 0.93	78.17 \pm 0.89
D_t	BEFORE	93.13 \pm 0.00					
	NegGrad	84.82 \pm 7.64	66.76 \pm 9.92	40.66 \pm 6.60	21.56 \pm 3.65	8.08 \pm 1.33	4.72 \pm 0.84
	RandL	84.33 \pm 2.72	82.71 \pm 1.18	75.24 \pm 4.03	59.91 \pm 2.60	43.35 \pm 2.96	28.88 \pm 1.21
	Bad Teacher	74.16 \pm 3.86	76.79 \pm 1.39	77.12 \pm 0.87	74.00 \pm 1.40	64.28 \pm 4.14	43.15 \pm 5.49
	Amnesiac	74.77 \pm 1.19	75.50 \pm 1.03	74.40 \pm 1.48	70.59 \pm 1.00	57.57 \pm 1.35	34.47 \pm 2.83
	δ -Targeted ¹	90.17 \pm 1.21	90.72 \pm 0.75	88.19 \pm 2.00	83.61 \pm 2.05	71.06 \pm 2.31	50.03 \pm 1.96
	δ -Targeted ³	88.06 \pm 1.55	89.64 \pm 0.74	89.30 \pm 0.63	88.17 \pm 0.62	82.68 \pm 1.00	69.63 \pm 1.05

Table 6: Means and standard deviations (percentage) of classification accuracy on CIFAR10 with ResNet-50.

Method	$ D_f =16$	$ D_f =32$	$ D_f =64$	$ D_f =128$	$ D_f =256$	$ D_f =512$
BEFORE	100.00 \pm 0.00					
NegGrad	0.00 \pm 0.00					
RandL	0.00 \pm 0.00					
D _f Bad Teacher	13.75 \pm 10.00	10.62 \pm 4.24	10.62 \pm 3.19	11.72 \pm 2.92	10.70 \pm 1.09	10.78 \pm 1.71
Amnesiac	0.00 \pm 0.00					
δ -Targeted ¹	0.00 \pm 0.00					
δ -Targeted ³	0.00 \pm 0.00					
BEFORE	100.00 \pm 0.00					
NegGrad	76.01 \pm 4.78	64.18 \pm 3.28	38.01 \pm 6.06	17.06 \pm 4.24	8.59 \pm 1.43	4.94 \pm 0.98
RandL	92.51 \pm 1.67	89.36 \pm 2.70	79.43 \pm 3.25	64.72 \pm 4.14	47.95 \pm 4.95	33.31 \pm 2.80
D _r Bad Teacher	70.56 \pm 2.01	68.54 \pm 3.39	65.62 \pm 1.04	54.88 \pm 4.27	36.45 \pm 3.05	23.18 \pm 4.47
Amnesiac	70.81 \pm 0.30	70.05 \pm 1.37	69.05 \pm 1.63	64.36 \pm 2.30	43.60 \pm 1.78	26.51 \pm 0.88
δ -Targeted ¹	98.92 \pm 0.44	98.55 \pm 0.51	96.54 \pm 0.73	90.86 \pm 1.42	80.69 \pm 3.27	58.80 \pm 2.02
δ -Targeted ³	96.31 \pm 0.74	96.96 \pm 0.95	97.23 \pm 0.94	96.36 \pm 0.91	92.09 \pm 0.78	75.74 \pm 0.82
BEFORE	94.02 \pm 0.00					
NegGrad	70.31 \pm 4.60	59.22 \pm 3.17	34.67 \pm 5.74	15.26 \pm 3.84	8.03 \pm 1.40	4.53 \pm 0.81
RandL	85.60 \pm 1.69	82.03 \pm 2.49	72.26 \pm 3.25	58.06 \pm 3.52	42.72 \pm 4.09	29.52 \pm 2.69
D _t Bad Teacher	68.93 \pm 2.02	66.55 \pm 3.42	64.12 \pm 1.23	53.19 \pm 4.25	35.59 \pm 3.11	22.50 \pm 4.55
Amnesiac	69.17 \pm 0.40	68.22 \pm 1.61	66.92 \pm 1.61	62.33 \pm 1.88	42.08 \pm 2.14	25.42 \pm 1.56
δ -Targeted ¹	91.55 \pm 0.39	91.15 \pm 0.42	88.92 \pm 0.74	83.04 \pm 1.44	72.93 \pm 3.24	52.51 \pm 2.22
δ -Targeted ³	89.04 \pm 0.61	89.75 \pm 0.78	89.90 \pm 0.91	88.90 \pm 0.98	84.37 \pm 1.01	68.58 \pm 0.88

Table 7: Means and standard deviations (percentage) of classification accuracy on CIFAR10 with ViT.

Method	$ D_f =16$	$ D_f =32$	$ D_f =64$	$ D_f =128$	$ D_f =256$	$ D_f =512$
BEFORE	93.75 \pm 6.85	94.38 \pm 3.64	95.00 \pm 3.62	95.16 \pm 2.59	94.77 \pm 2.27	93.48 \pm 1.33
NegGrad	8.75 \pm 10.90	10.00 \pm 2.34	8.75 \pm 2.12	9.38 \pm 2.21	10.86 \pm 1.43	10.43 \pm 1.41
RandL	0.00 \pm 0.00					
D _f Bad Teacher	12.50 \pm 5.59	6.88 \pm 4.59	12.19 \pm 4.57	15.78 \pm 5.49	18.59 \pm 2.86	25.94 \pm 2.20
Amnesiac	0.00 \pm 0.00					
δ -Targeted ¹	0.00 \pm 0.00					
δ -Targeted ³	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.55 \pm 0.40	2.70 \pm 1.01
BEFORE	93.64 \pm 0.00	93.64 \pm 0.01	93.64 \pm 0.01	93.64 \pm 0.01	93.63 \pm 0.01	93.64 \pm 0.01
NegGrad	40.34 \pm 8.32	21.06 \pm 6.47	14.10 \pm 7.21	10.05 \pm 0.10	9.99 \pm 0.01	9.99 \pm 0.01
RandL	79.26 \pm 1.48	74.43 \pm 1.08	61.48 \pm 3.34	47.20 \pm 4.53	29.39 \pm 1.79	18.96 \pm 2.11
D _r Bad Teacher	85.98 \pm 0.44	84.69 \pm 0.55	80.46 \pm 1.07	75.41 \pm 0.97	65.98 \pm 1.38	53.28 \pm 3.38
Amnesiac	81.04 \pm 0.32	78.16 \pm 0.53	74.15 \pm 1.13	68.07 \pm 1.19	58.34 \pm 0.83	40.50 \pm 2.02
δ -Targeted ¹	85.74 \pm 3.65	83.81 \pm 2.10	74.86 \pm 3.57	70.53 \pm 3.32	65.69 \pm 2.58	55.62 \pm 2.66
δ -Targeted ³	87.21 \pm 1.78	85.96 \pm 1.10	81.13 \pm 2.18	78.66 \pm 1.93	75.13 \pm 2.29	67.42 \pm 1.52
BEFORE	81.69 \pm 0.00					
NegGrad	38.11 \pm 7.29	20.33 \pm 5.93	13.80 \pm 6.51	10.03 \pm 0.11	10.00 \pm 0.12	9.91 \pm 0.16
RandL	70.53 \pm 1.21	65.85 \pm 0.86	54.26 \pm 3.02	42.25 \pm 3.84	26.26 \pm 1.69	16.90 \pm 1.63
D _t Bad Teacher	77.51 \pm 0.52	76.52 \pm 0.45	72.53 \pm 0.82	68.53 \pm 0.59	60.51 \pm 1.21	50.43 \pm 2.82
Amnesiac	72.48 \pm 0.30	69.71 \pm 0.48	66.16 \pm 1.09	60.93 \pm 0.87	52.44 \pm 0.83	36.84 \pm 1.96
δ -Targeted ¹	75.47 \pm 2.67	73.91 \pm 1.47	66.41 \pm 2.81	62.63 \pm 2.84	58.08 \pm 2.25	49.50 \pm 2.40
δ -Targeted ³	76.68 \pm 1.25	75.84 \pm 0.81	71.66 \pm 1.76	69.43 \pm 1.44	66.20 \pm 1.97	59.58 \pm 1.25

1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508

1509 **B.2 Performance on CIFAR100 with Different Neural Network Structures** 1567

1510

1511

1512 **Table 8: Means and standard deviations (percentage) of classification accuracy on CIFAR100 with ResNet-18.** 1570

1513

Method	$ D_f =16$	$ D_f =32$	$ D_f =64$	$ D_f =128$	$ D_f =256$	$ D_f =512$	
D_f	<i>BEFORE</i>	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	99.84 \pm 0.31	99.92 \pm 0.16	99.96 \pm 0.08
	NegGrad	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	20.00 \pm 40.00	0.00 \pm 0.00	0.00 \pm 0.00
	RandL	0.00 \pm 0.00					
	Bad Teacher	0.00 \pm 0.00	0.62 \pm 1.25	0.94 \pm 1.25	0.47 \pm 0.62	1.33 \pm 0.19	0.86 \pm 0.23
	Amnesiac	0.00 \pm 0.00					
	δ -Targeted ¹	0.00 \pm 0.00					
	δ -Targeted ³	0.00 \pm 0.00					
D_r	<i>BEFORE</i>	99.96 \pm 0.00					
	NegGrad	83.71 \pm 3.40	86.98 \pm 0.84	84.97 \pm 1.28	83.83 \pm 8.28	70.57 \pm 1.72	55.92 \pm 0.95
	RandL	91.86 \pm 2.90	95.32 \pm 0.86	95.40 \pm 0.57	94.58 \pm 1.13	92.86 \pm 0.40	89.79 \pm 0.54
	Bad Teacher	36.13 \pm 1.50	36.75 \pm 1.51	35.45 \pm 1.66	29.34 \pm 2.15	25.57 \pm 2.45	16.33 \pm 2.24
	Amnesiac	35.47 \pm 1.62	34.69 \pm 0.65	33.22 \pm 2.57	31.17 \pm 0.89	22.39 \pm 0.88	11.80 \pm 0.63
	δ -Targeted ¹	98.17 \pm 0.90	98.99 \pm 0.35	98.84 \pm 0.10	98.16 \pm 0.71	97.34 \pm 0.37	94.53 \pm 0.56
	δ -Targeted ³	94.89 \pm 2.10	97.10 \pm 0.69	97.05 \pm 0.46	96.49 \pm 1.05	96.06 \pm 0.20	95.36 \pm 0.41
D_t	<i>BEFORE</i>	71.94 \pm 0.00					
	NegGrad	56.30 \pm 2.06	57.75 \pm 0.63	56.30 \pm 1.00	56.35 \pm 7.57	46.14 \pm 1.10	36.38 \pm 0.44
	RandL	60.95 \pm 2.05	64.02 \pm 0.57	64.28 \pm 0.36	63.21 \pm 0.71	60.74 \pm 0.46	56.39 \pm 0.23
	Bad Teacher	30.66 \pm 0.90	31.33 \pm 1.25	30.04 \pm 1.35	24.86 \pm 1.82	21.32 \pm 1.97	13.97 \pm 1.77
	Amnesiac	30.54 \pm 1.54	29.85 \pm 0.98	28.35 \pm 2.16	26.93 \pm 0.94	19.26 \pm 0.96	10.08 \pm 0.42
	δ -Targeted ¹	66.27 \pm 1.38	68.06 \pm 0.52	68.08 \pm 0.16	67.36 \pm 0.38	66.08 \pm 0.30	62.87 \pm 0.71
	δ -Targeted ³	63.38 \pm 1.72	65.89 \pm 0.60	66.52 \pm 0.35	66.22 \pm 0.60	65.45 \pm 0.23	64.06 \pm 0.55

1536

1537

1538

1539

1540 **Table 9: Means and standard deviations (percentage) of classification accuracy on CIFAR100 with ResNet-50.** 1598

1541

Method	$ D_f =16$	$ D_f =32$	$ D_f =64$	$ D_f =128$	$ D_f =256$	$ D_f =512$	
D_f	<i>BEFORE</i>	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	99.84 \pm 0.31	99.92 \pm 0.16	99.96 \pm 0.08
	NegGrad	0.00 \pm 0.00					
	RandL	0.00 \pm 0.00					
	Bad Teacher	0.00 \pm 0.00	0.62 \pm 1.25	1.25 \pm 1.17	1.25 \pm 0.80	1.33 \pm 0.19	1.41 \pm 0.19
	Amnesiac	0.00 \pm 0.00					
	δ -Targeted ¹	0.00 \pm 0.00					
	δ -Targeted ³	0.00 \pm 0.00					
D_r	<i>BEFORE</i>	99.97 \pm 0.00					
	NegGrad	66.36 \pm 7.53	66.55 \pm 5.26	60.86 \pm 3.08	57.38 \pm 3.42	51.95 \pm 3.63	43.16 \pm 1.36
	RandL	87.54 \pm 6.39	91.20 \pm 2.52	92.83 \pm 1.95	93.29 \pm 1.47	92.44 \pm 0.53	89.08 \pm 0.85
	Bad Teacher	23.47 \pm 1.07	23.83 \pm 2.11	22.25 \pm 0.71	11.25 \pm 1.07	5.22 \pm 0.82	3.36 \pm 0.51
	Amnesiac	24.92 \pm 0.98	25.12 \pm 1.31	24.37 \pm 0.62	21.51 \pm 1.67	12.26 \pm 1.60	7.36 \pm 0.56
	δ -Targeted ¹	98.36 \pm 1.39	97.07 \pm 1.26	97.27 \pm 1.06	97.36 \pm 0.66	96.87 \pm 0.63	93.58 \pm 0.74
	δ -Targeted ³	93.71 \pm 3.49	94.68 \pm 1.86	95.37 \pm 1.43	95.55 \pm 1.36	95.46 \pm 0.53	94.27 \pm 0.67
D_t	<i>BEFORE</i>	74.67 \pm 0.00					
	NegGrad	45.94 \pm 4.63	45.44 \pm 3.18	41.33 \pm 2.02	38.76 \pm 2.07	35.34 \pm 2.76	29.79 \pm 1.02
	RandL	60.14 \pm 4.44	62.62 \pm 2.06	64.22 \pm 1.76	64.39 \pm 1.40	62.79 \pm 0.48	58.55 \pm 0.41
	Bad Teacher	21.45 \pm 0.90	21.28 \pm 2.17	20.22 \pm 0.75	9.70 \pm 1.08	4.41 \pm 0.74	2.62 \pm 0.56
	Amnesiac	22.30 \pm 0.69	22.33 \pm 1.07	21.98 \pm 0.69	19.08 \pm 1.49	10.91 \pm 1.58	6.56 \pm 0.77
	δ -Targeted ¹	69.77 \pm 2.57	67.91 \pm 1.18	68.89 \pm 1.01	69.18 \pm 0.87	68.55 \pm 0.67	64.65 \pm 0.49
	δ -Targeted ³	64.80 \pm 3.68	65.76 \pm 1.73	67.09 \pm 1.40	67.48 \pm 1.47	67.26 \pm 0.43	65.51 \pm 0.60

1563

1564

1565

1566

Table 10: Means and standard deviations (percentage) of classification accuracy on CIFAR100 with ViT.

Method	$ D_f =16$	$ D_f =32$	$ D_f =64$	$ D_f =128$	$ D_f =256$	$ D_f =512$	
D_f	<i>BEFORE</i>	86.25 \pm 8.29	81.88 \pm 8.93	83.44 \pm 2.54	84.84 \pm 2.19	84.30 \pm 0.62	84.84 \pm 1.35
	NegGrad	1.25 \pm 2.50	0.62 \pm 1.25	0.62 \pm 0.77	1.09 \pm 0.38	0.78 \pm 0.49	0.70 \pm 0.47
	RandL	0.00 \pm 0.00					
	Bad Teacher	3.75 \pm 5.00	2.50 \pm 2.34	4.06 \pm 1.88	5.00 \pm 1.45	6.72 \pm 0.76	10.35 \pm 0.54
	Amnesiac	0.00 \pm 0.00					
	δ -Targeted ¹	0.00 \pm 0.00					
	δ -Targeted ³	0.00 \pm 0.00					
D_r	<i>BEFORE</i>	85.18 \pm 0.01	85.18 \pm 0.01	85.19 \pm 0.02	85.18 \pm 0.00	85.19 \pm 0.00	85.19 \pm 0.01
	NegGrad	38.64 \pm 2.76	25.91 \pm 1.68	13.58 \pm 3.10	1.55 \pm 1.09	1.00 \pm 0.00	1.00 \pm 0.01
	RandL	66.80 \pm 3.46	63.40 \pm 2.19	56.68 \pm 1.36	47.10 \pm 2.25	36.13 \pm 2.18	22.14 \pm 0.82
	Bad Teacher	74.08 \pm 0.38	71.23 \pm 0.60	63.16 \pm 1.72	52.82 \pm 0.80	39.30 \pm 1.25	26.99 \pm 1.19
	Amnesiac	64.19 \pm 1.46	62.44 \pm 1.16	57.94 \pm 1.15	53.96 \pm 0.99	46.12 \pm 0.75	31.17 \pm 1.18
	δ -Targeted ¹	83.65 \pm 0.52	82.49 \pm 1.10	79.90 \pm 0.99	76.41 \pm 1.00	73.74 \pm 1.49	70.48 \pm 1.16
	δ -Targeted ³	84.04 \pm 0.74	82.97 \pm 1.15	80.86 \pm 0.71	78.79 \pm 0.71	77.57 \pm 1.36	75.54 \pm 0.65
D_t	<i>BEFORE</i>	54.59 \pm 0.00					
	NegGrad	28.75 \pm 1.70	20.25 \pm 1.36	11.22 \pm 2.61	1.45 \pm 0.86	1.01 \pm 0.06	0.98 \pm 0.00
	RandL	45.19 \pm 1.59	42.82 \pm 1.23	38.65 \pm 1.03	32.59 \pm 1.49	25.17 \pm 1.46	15.58 \pm 0.43
	Bad Teacher	50.66 \pm 0.48	49.06 \pm 0.30	44.50 \pm 0.66	38.84 \pm 0.74	30.35 \pm 1.05	22.14 \pm 0.85
	Amnesiac	43.38 \pm 0.59	42.34 \pm 0.61	39.70 \pm 0.59	37.25 \pm 0.78	31.76 \pm 0.50	22.45 \pm 0.57
	δ -Targeted ¹	53.69 \pm 0.20	53.25 \pm 0.49	51.63 \pm 0.59	49.86 \pm 0.41	47.92 \pm 0.94	46.00 \pm 0.78
	δ -Targeted ³	53.92 \pm 0.24	53.37 \pm 0.36	52.14 \pm 0.48	51.31 \pm 0.32	50.21 \pm 0.70	48.73 \pm 0.70

B.3 Performance on Fashion-MNIST with Different Neural Network Structures**Table 11: Means and standard deviations (percentage) of classification accuracy on Fashion-MNIST with ResNet-18.**

Method	$ D_f =16$	$ D_f =32$	$ D_f =64$	$ D_f =128$	$ D_f =256$	$ D_f =512$	
D_f	<i>BEFORE</i>	95.00 \pm 4.68	94.38 \pm 3.64	94.69 \pm 1.59	95.47 \pm 0.91	96.02 \pm 0.67	96.48 \pm 0.86
	NegGrad	0.00 \pm 0.00					
	RandL	0.00 \pm 0.00					
	Bad Teacher	12.50 \pm 6.85	6.88 \pm 2.34	8.44 \pm 4.38	9.84 \pm 3.51	12.89 \pm 3.62	13.95 \pm 5.09
	Amnesiac	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.94 \pm 1.37	3.24 \pm 1.88
	δ -Targeted ¹	0.00 \pm 0.00					
	δ -Targeted ³	0.00 \pm 0.00					
D_r	<i>BEFORE</i>	96.90 \pm 0.00	96.90 \pm 0.00	96.90 \pm 0.01	96.90 \pm 0.00	96.90 \pm 0.00	96.90 \pm 0.01
	NegGrad	73.65 \pm 1.36	61.92 \pm 3.62	45.79 \pm 3.52	26.79 \pm 0.93	16.05 \pm 1.11	8.77 \pm 1.49
	RandL	86.51 \pm 1.35	85.38 \pm 1.46	75.68 \pm 1.65	60.56 \pm 1.10	43.32 \pm 1.95	28.76 \pm 1.22
	Bad Teacher	86.49 \pm 0.40	86.07 \pm 0.29	85.71 \pm 0.52	80.54 \pm 2.86	68.33 \pm 2.65	55.75 \pm 3.67
	Amnesiac	85.68 \pm 0.40	84.32 \pm 0.77	81.41 \pm 1.03	76.38 \pm 1.07	59.85 \pm 3.75	42.45 \pm 3.57
	δ -Targeted ¹	93.94 \pm 0.50	92.13 \pm 0.68	87.06 \pm 0.27	76.28 \pm 1.39	59.01 \pm 3.18	39.59 \pm 1.16
	δ -Targeted ³	92.93 \pm 0.98	93.35 \pm 0.41	92.40 \pm 0.57	87.61 \pm 1.21	76.71 \pm 1.94	58.28 \pm 2.43
D_t	<i>BEFORE</i>	91.65 \pm 0.00					
	NegGrad	69.50 \pm 1.73	58.26 \pm 4.21	42.21 \pm 3.13	23.90 \pm 0.98	13.82 \pm 1.02	7.57 \pm 1.16
	RandL	82.03 \pm 1.07	81.17 \pm 1.37	71.07 \pm 2.18	56.25 \pm 1.06	39.02 \pm 2.11	24.99 \pm 1.30
	Bad Teacher	85.18 \pm 0.34	84.56 \pm 0.36	84.27 \pm 0.40	78.98 \pm 2.81	67.08 \pm 2.44	55.04 \pm 3.51
	Amnesiac	84.37 \pm 0.56	83.07 \pm 0.54	80.05 \pm 1.17	74.67 \pm 0.90	58.42 \pm 3.83	41.74 \pm 3.85
	δ -Targeted ¹	88.91 \pm 0.48	87.35 \pm 0.70	82.36 \pm 0.48	71.60 \pm 2.03	54.48 \pm 3.36	35.16 \pm 1.24
	δ -Targeted ³	88.12 \pm 0.79	88.57 \pm 0.37	87.42 \pm 0.58	82.74 \pm 1.09	72.01 \pm 1.92	53.66 \pm 2.49

1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740

Table 12: Means and standard deviations (percentage) of classification accuracy on Fashion-MNIST with ResNet-50.

	Method	$ D_f =16$	$ D_f =32$	$ D_f =64$	$ D_f =128$	$ D_f =256$	$ D_f =512$	
D_f	BEFORE	100.00 \pm 0.00	98.75 \pm 1.53	97.81 \pm 1.59	97.66 \pm 1.48	96.95 \pm 1.06	97.30 \pm 0.52	1799
	NegGrad	0.00 \pm 0.00	1800					
	RandL	0.00 \pm 0.00	1801					
	Bad Teacher	11.25 \pm 7.29	11.25 \pm 4.24	11.88 \pm 2.90	10.94 \pm 2.47	13.59 \pm 1.79	12.66 \pm 2.05	1802
	Amnesiac	0.00 \pm 0.00	0.00 \pm 0.00	0.94 \pm 1.88	0.47 \pm 0.94	6.25 \pm 4.98	3.05 \pm 1.36	1803
	δ -Targeted ¹	0.00 \pm 0.00	1804					
	δ -Targeted ³	0.00 \pm 0.00	1805					
D_r	BEFORE	97.10 \pm 0.00	1806					
	NegGrad	68.70 \pm 5.09	51.92 \pm 6.19	34.08 \pm 3.29	18.40 \pm 1.46	9.01 \pm 0.71	5.29 \pm 0.41	1807
	RandL	81.25 \pm 6.89	75.48 \pm 3.47	64.52 \pm 4.04	46.83 \pm 2.00	29.98 \pm 1.67	20.67 \pm 0.90	1808
	Bad Teacher	83.84 \pm 0.50	83.03 \pm 1.65	77.70 \pm 2.68	69.85 \pm 4.52	51.28 \pm 5.79	36.32 \pm 3.92	1809
	Amnesiac	83.74 \pm 1.30	82.68 \pm 0.98	78.56 \pm 1.80	73.40 \pm 0.83	60.89 \pm 2.75	39.38 \pm 1.59	1810
	δ -Targeted ¹	94.11 \pm 0.63	88.65 \pm 2.09	80.80 \pm 2.20	68.19 \pm 3.23	51.23 \pm 3.56	34.98 \pm 1.17	1811
	δ -Targeted ³	93.96 \pm 1.11	93.56 \pm 0.34	90.87 \pm 1.12	83.50 \pm 2.48	69.57 \pm 1.73	49.50 \pm 0.85	1812
D_t	BEFORE	91.79 \pm 0.00	1813					
	NegGrad	64.50 \pm 5.13	47.82 \pm 6.14	30.65 \pm 3.27	15.94 \pm 1.63	7.76 \pm 0.61	4.53 \pm 0.26	1814
	RandL	76.48 \pm 6.88	70.32 \pm 3.42	59.76 \pm 4.00	42.36 \pm 1.81	26.24 \pm 1.62	17.54 \pm 1.10	1815
	Bad Teacher	82.79 \pm 0.52	81.58 \pm 1.94	76.59 \pm 2.54	69.00 \pm 4.55	50.54 \pm 5.64	35.94 \pm 3.93	1816
	Amnesiac	82.81 \pm 1.29	81.62 \pm 1.12	77.85 \pm 1.79	72.95 \pm 0.79	60.62 \pm 2.94	38.86 \pm 1.98	1817
	δ -Targeted ¹	89.19 \pm 0.53	83.47 \pm 2.55	75.66 \pm 2.41	62.94 \pm 3.07	46.41 \pm 3.53	31.12 \pm 0.94	1818
	δ -Targeted ³	89.44 \pm 0.73	88.87 \pm 0.35	86.05 \pm 1.28	78.45 \pm 2.57	64.35 \pm 1.61	45.15 \pm 0.77	1819

Table 13: Means and standard deviations (percentage) of classification accuracy on Fashion-MNIST with ViT.

	Method	$ D_f =16$	$ D_f =32$	$ D_f =64$	$ D_f =128$	$ D_f =256$	$ D_f =512$	
D_f	BEFORE	93.75 \pm 5.59	93.12 \pm 4.59	91.88 \pm 2.69	91.72 \pm 1.17	92.42 \pm 0.58	93.48 \pm 0.64	1820
	NegGrad	8.75 \pm 9.35	8.12 \pm 9.19	10.94 \pm 6.63	6.72 \pm 5.13	10.08 \pm 2.87	10.43 \pm 1.51	1821
	RandL	0.00 \pm 0.00	1.37 \pm 0.35	1822				
	Bad Teacher	0.00 \pm 0.00	1823					
	Amnesiac	0.00 \pm 0.00	1.25 \pm 2.50	6.56 \pm 2.50	18.44 \pm 3.37	29.92 \pm 2.82	25.74 \pm 2.58	1824
	δ -Targeted ¹	0.00 \pm 0.00	0.00 \pm 0.00	1.88 \pm 1.17	3.28 \pm 1.94	13.67 \pm 3.91	52.50 \pm 4.18	1825
	δ -Targeted ³	33.75 \pm 8.48	31.25 \pm 6.56	46.88 \pm 6.01	67.19 \pm 0.49	76.64 \pm 2.35	77.54 \pm 2.08	1826
D_r	BEFORE	94.11 \pm 0.00	94.11 \pm 0.00	94.11 \pm 0.01	94.11 \pm 0.00	94.11 \pm 0.00	94.11 \pm 0.01	1827
	NegGrad	19.60 \pm 7.44	9.94 \pm 1.37	11.50 \pm 4.18	8.39 \pm 3.88	10.00 \pm 0.01	10.00 \pm 0.01	1828
	RandL	48.39 \pm 6.41	31.38 \pm 5.46	20.56 \pm 4.46	13.27 \pm 2.49	8.74 \pm 2.35	7.19 \pm 1.05	1829
	Bad Teacher	67.16 \pm 9.11	59.51 \pm 4.55	50.04 \pm 3.26	39.13 \pm 2.33	27.45 \pm 1.21	21.05 \pm 0.66	1830
	Amnesiac	75.10 \pm 6.75	73.45 \pm 3.09	64.50 \pm 2.36	57.24 \pm 1.96	48.30 \pm 1.54	38.16 \pm 1.86	1831
	δ -Targeted ¹	87.21 \pm 0.54	84.90 \pm 1.33	80.34 \pm 1.15	73.08 \pm 1.23	62.07 \pm 1.20	61.94 \pm 2.66	1832
	δ -Targeted ³	90.53 \pm 0.31	89.02 \pm 0.67	86.85 \pm 0.67	84.00 \pm 0.38	80.82 \pm 1.02	77.31 \pm 1.44	1833
D_t	BEFORE	92.41 \pm 0.00	1834					
	NegGrad	19.58 \pm 7.47	10.03 \pm 1.46	11.58 \pm 4.10	8.51 \pm 3.93	10.15 \pm 0.12	10.06 \pm 0.19	1835
	RandL	47.47 \pm 6.61	30.98 \pm 5.32	20.20 \pm 4.20	13.16 \pm 2.61	8.71 \pm 2.08	6.87 \pm 0.96	1836
	Bad Teacher	65.85 \pm 9.02	58.35 \pm 4.44	48.81 \pm 3.24	37.86 \pm 2.46	26.34 \pm 1.54	20.44 \pm 0.78	1837
	Amnesiac	73.54 \pm 6.76	71.91 \pm 3.05	62.96 \pm 2.56	55.70 \pm 2.01	46.49 \pm 1.51	36.23 \pm 1.94	1838
	δ -Targeted ¹	85.87 \pm 0.77	83.72 \pm 1.38	79.11 \pm 1.22	71.81 \pm 1.36	61.22 \pm 1.35	61.48 \pm 2.44	1839
	δ -Targeted ³	89.09 \pm 0.31	87.58 \pm 0.71	85.36 \pm 0.81	82.90 \pm 0.57	80.05 \pm 1.07	76.80 \pm 1.38	1840

B.4 Performance on Tiny-ImageNet with Different Neural Network Structures

Table 14: Means and standard deviations (percentage) of classification accuracy on Tiny-ImageNet with ResNet-18.

Method	$ D_f =16$	$ D_f =32$	$ D_f =64$	$ D_f =128$	$ D_f =256$	$ D_f =512$
D_f	<i>BEFORE</i>	85.00 ± 6.37	85.62 ± 3.75	84.69 ± 2.07	84.84 ± 1.06	84.06 ± 1.17
	NegGrad	0.00 ± 0.00				
	RandL	0.00 ± 0.00				
	Bad Teacher	1.25 ± 2.50	0.62 ± 1.25	0.00 ± 0.00	0.62 ± 0.58	0.31 ± 0.46
	Amnesiac	0.00 ± 0.00				
	δ -Targeted ¹	0.00 ± 0.00				
	δ -Targeted ³	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.04 ± 0.08
D_r	<i>BEFORE</i>	85.01 ± 0.01	85.01 ± 0.00	85.01 ± 0.00	85.02 ± 0.00	85.02 ± 0.00
	NegGrad	58.03 ± 1.68	61.53 ± 1.86	63.25 ± 0.68	58.79 ± 0.71	53.26 ± 0.56
	RandL	62.68 ± 1.90	69.22 ± 2.68	74.88 ± 1.21	76.67 ± 0.44	76.29 ± 0.48
	Bad Teacher	18.82 ± 0.57	19.46 ± 0.65	19.73 ± 1.07	18.01 ± 0.79	12.37 ± 1.02
	Amnesiac	12.62 ± 0.64	12.67 ± 0.54	12.68 ± 0.78	12.18 ± 0.57	11.63 ± 0.54
	δ -Targeted ¹	83.02 ± 1.07	83.43 ± 0.52	83.60 ± 0.38	83.36 ± 0.26	82.93 ± 0.18
	δ -Targeted ³	76.88 ± 4.04	77.96 ± 2.31	80.45 ± 0.71	81.89 ± 0.33	82.77 ± 0.22
D_t	<i>BEFORE</i>	52.71 ± 0.00				
	NegGrad	40.47 ± 1.05	42.26 ± 0.87	42.37 ± 0.81	39.40 ± 0.77	35.90 ± 0.42
	RandL	43.58 ± 1.17	46.22 ± 1.31	48.25 ± 0.42	48.39 ± 0.40	47.61 ± 0.49
	Bad Teacher	17.81 ± 0.62	18.90 ± 0.51	18.87 ± 1.42	17.08 ± 1.00	11.72 ± 0.97
	Amnesiac	12.22 ± 0.46	12.15 ± 0.63	12.07 ± 0.72	11.40 ± 0.42	10.65 ± 0.82
	δ -Targeted ¹	52.00 ± 0.34	52.18 ± 0.23	52.04 ± 0.24	51.75 ± 0.30	51.57 ± 0.25
	δ -Targeted ³	49.63 ± 1.32	50.08 ± 0.64	50.97 ± 0.49	51.24 ± 0.41	51.40 ± 0.40

Table 15: Means and standard deviations (percentage) of classification accuracy on Tiny-ImageNet with ResNet-50.

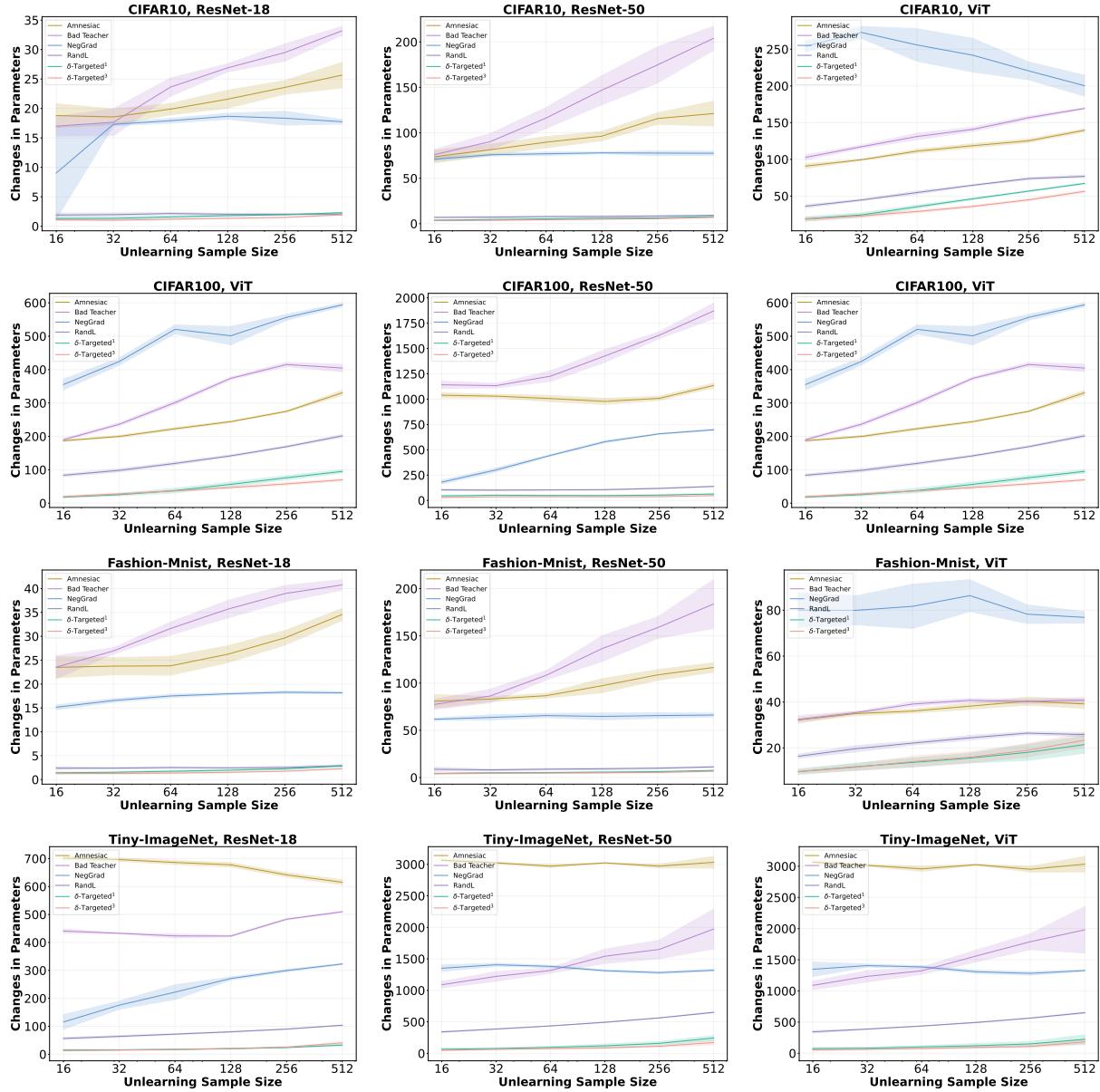
Method	$ D_f =16$	$ D_f =32$	$ D_f =64$	$ D_f =128$	$ D_f =256$	$ D_f =512$
D_f	<i>BEFORE</i>	28.12 ± 9.38	29.69 ± 4.69	28.91 ± 0.78	30.86 ± 5.08	32.23 ± 3.71
	NegGrad	0.00 ± 0.00				
	RandL	0.00 ± 0.00				
	Bad Teacher	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.31 ± 0.38	0.31 ± 0.29
	Amnesiac	0.00 ± 0.00				
	δ -Targeted ¹	0.00 ± 0.00				
	δ -Targeted ³	0.00 ± 0.00				
D_r	<i>BEFORE</i>	53.04 ± 0.02	53.05 ± 0.00	53.05 ± 0.00	53.06 ± 0.00	53.05 ± 0.00
	NegGrad	38.70 ± 1.30	42.01 ± 0.85	40.98 ± 0.73	36.26 ± 1.20	30.46 ± 0.75
	RandL	39.28 ± 0.61	38.70 ± 0.88	39.13 ± 0.67	40.85 ± 0.63	42.73 ± 0.38
	Bad Teacher	28.46 ± 0.77	27.44 ± 2.26	28.75 ± 0.56	26.27 ± 0.32	21.78 ± 0.28
	Amnesiac	29.35 ± 0.61	29.50 ± 0.46	29.59 ± 0.57	29.07 ± 0.39	27.57 ± 0.35
	δ -Targeted ¹	51.35 ± 0.68	50.34 ± 0.60	49.27 ± 0.43	47.84 ± 1.11	46.69 ± 0.84
	δ -Targeted ³	51.65 ± 0.81	51.29 ± 0.33	50.74 ± 0.20	49.88 ± 0.47	48.84 ± 0.61
D_t	<i>BEFORE</i>	52.17 ± 0.00				
	NegGrad	37.86 ± 0.91	41.03 ± 0.91	40.44 ± 0.85	35.77 ± 1.19	29.85 ± 0.83
	RandL	38.51 ± 0.72	38.15 ± 0.69	38.67 ± 0.63	40.35 ± 0.34	41.80 ± 0.37
	Bad Teacher	28.60 ± 1.01	27.45 ± 2.33	28.70 ± 0.88	26.14 ± 0.50	21.80 ± 0.16
	Amnesiac	29.36 ± 0.65	29.38 ± 0.43	29.65 ± 0.64	28.93 ± 0.43	27.40 ± 0.37
	δ -Targeted ¹	50.54 ± 0.80	49.73 ± 0.68	48.46 ± 0.34	46.88 ± 1.11	45.77 ± 0.78
	δ -Targeted ³	50.89 ± 0.90	50.68 ± 0.39	50.02 ± 0.19	48.97 ± 0.60	47.86 ± 0.63

1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972

Table 16: Means and standard deviations (percentage) of classification accuracy on Tiny-ImageNet with ViT.

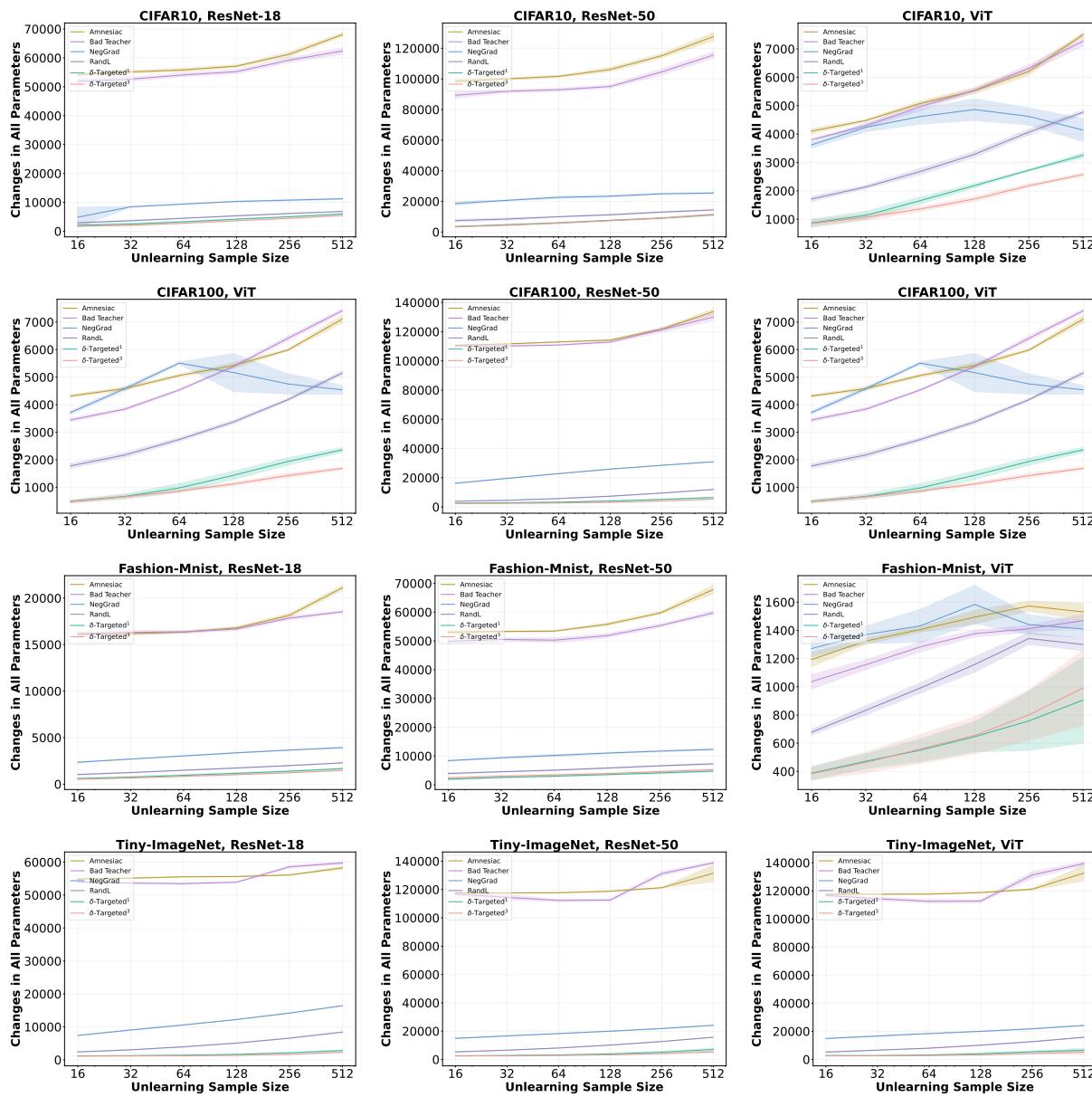
Method	$ D_f =16$	$ D_f =32$	$ D_f =64$	$ D_f =128$	$ D_f =256$	$ D_f =512$
<i>BEFORE</i>	49.75 \pm 1.25	49.50 \pm 0.50	50.05 \pm 2.25	50.10 \pm 0.40	51.55 \pm 1.45	47.00 \pm 0.00
NegGrad	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00
<i>D_f</i>	RandL	0.00 \pm 0.00				
<i>D_f</i>	Bad Teacher	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.31 \pm 0.38	0.31 \pm 0.29
	Amnesiac	0.00 \pm 0.00				
	δ -Targeted ¹	0.00 \pm 0.00				
	δ -Targeted ³	0.00 \pm 0.00				
<i>BEFORE</i>	54.89 \pm 0.00	54.89 \pm 0.00	54.89 \pm 0.00	54.89 \pm 0.00	54.89 \pm 0.00	54.89 \pm 0.00
NegGrad	37.60 \pm 1.96	43.18 \pm 5.69	44.17 \pm 3.90	36.36 \pm 6.14	30.96 \pm 4.92	29.29 \pm 5.16
<i>D_r</i>	RandL	36.90 \pm 4.11	35.76 \pm 3.72	37.21 \pm 4.06	40.06 \pm 7.68	43.20 \pm 4.56
<i>D_r</i>	Bad Teacher	28.36 \pm 2.01	29.74 \pm 3.41	25.26 \pm 1.75	25.49 \pm 6.03	23.50 \pm 3.93
	Amnesiac	29.61 \pm 3.36	27.86 \pm 3.56	26.82 \pm 2.44	27.45 \pm 5.00	28.32 \pm 2.86
	δ -Targeted ¹	50.67 \pm 2.84	49.30 \pm 5.14	45.94 \pm 3.66	46.37 \pm 2.44	49.24 \pm 3.04
	δ -Targeted ³	54.00 \pm 7.63	49.75 \pm 2.40	49.18 \pm 4.68	48.39 \pm 4.99	48.51 \pm 6.30
<i>BEFORE</i>	52.34 \pm 0.00	52.34 \pm 0.00	52.34 \pm 0.00	52.34 \pm 0.00	52.34 \pm 0.00	52.34 \pm 0.00
NegGrad	39.64 \pm 4.54	37.96 \pm 2.64	37.58 \pm 2.14	39.75 \pm 2.97	30.48 \pm 4.31	31.76 \pm 1.83
<i>D_t</i>	RandL	39.02 \pm 3.22	37.39 \pm 3.91	35.73 \pm 2.03	38.69 \pm 2.03	41.61 \pm 4.28
<i>D_t</i>	Bad Teacher	27.57 \pm 5.18	29.71 \pm 2.74	31.72 \pm 5.41	25.34 \pm 5.20	23.34 \pm 3.93
	Amnesiac	27.35 \pm 2.90	32.21 \pm 2.56	28.33 \pm 4.09	32.29 \pm 4.37	31.35 \pm 3.69
	δ -Targeted ¹	51.31 \pm 1.81	50.13 \pm 3.21	48.87 \pm 7.05	47.13 \pm 3.20	46.59 \pm 3.44
	δ -Targeted ³	49.60 \pm 2.88	50.42 \pm 1.82	48.52 \pm 5.52	52.09 \pm 3.99	48.85 \pm 2.92

2089 B.5 Changes in Last-Layer Parameters for Different Methods.



2134 **Figure 4:** The magnitude of parameter changes across the model, quantified by the L_1 -norm. We use $\Delta\theta$ to denote the change in
2135 all parameters and $\Delta\phi$ to denote the change in the last layer's parameters. Our methods lead to the smallest change to models'
2136 parameters on different datasets.

2205 B.6 Changes in All Parameters for Different Methods.



2250 **Figure 5:** The parameter changes across different neural network structures. The changes are quantified by the L_1 -norm. Our
2251 method leads to the smallest change to models' parameters on different datasets and cross different neural network structures.
2252

B.7 Relations between Changes in Last-Layer Parameters and Changes of All Parameters.

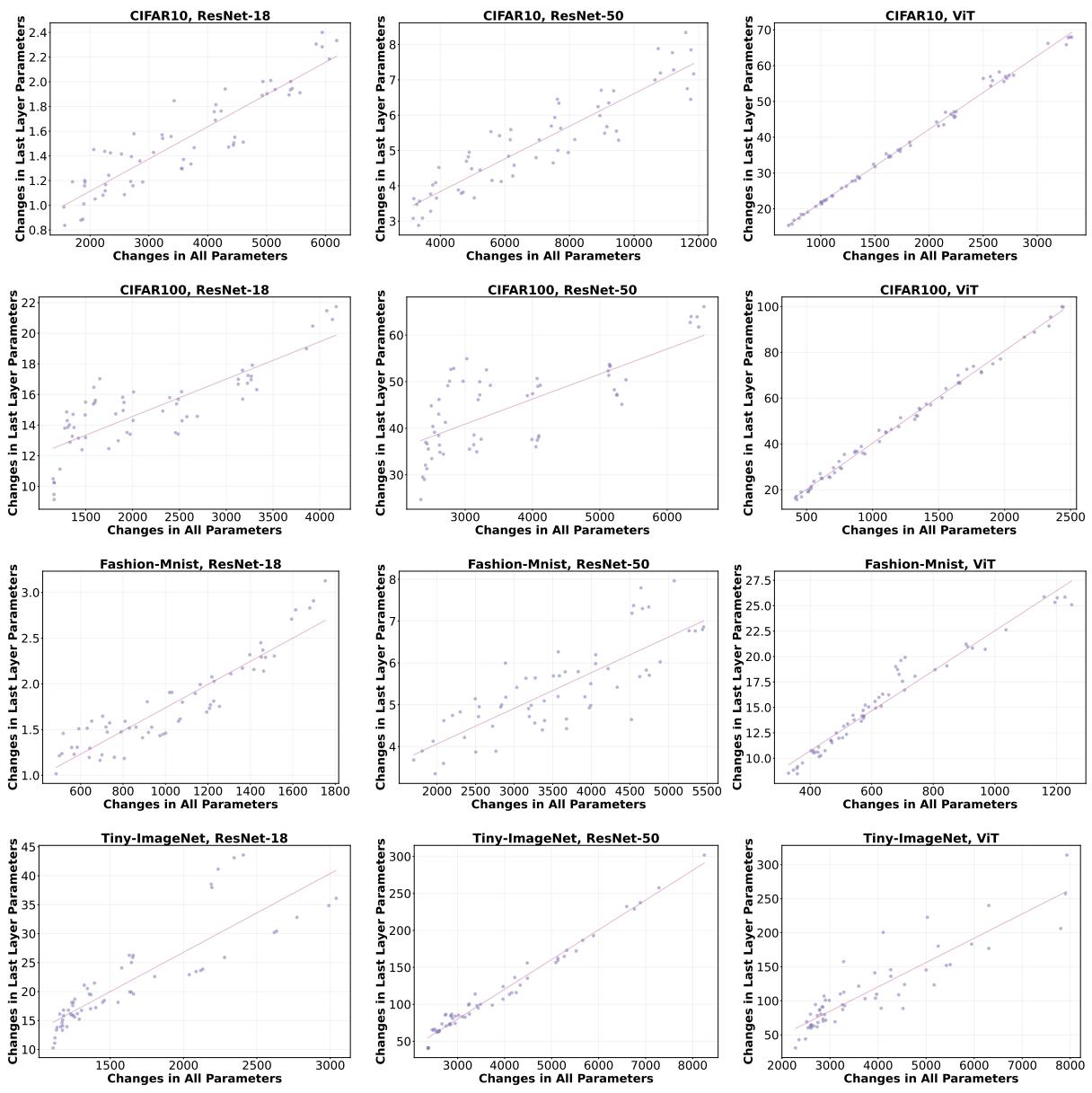
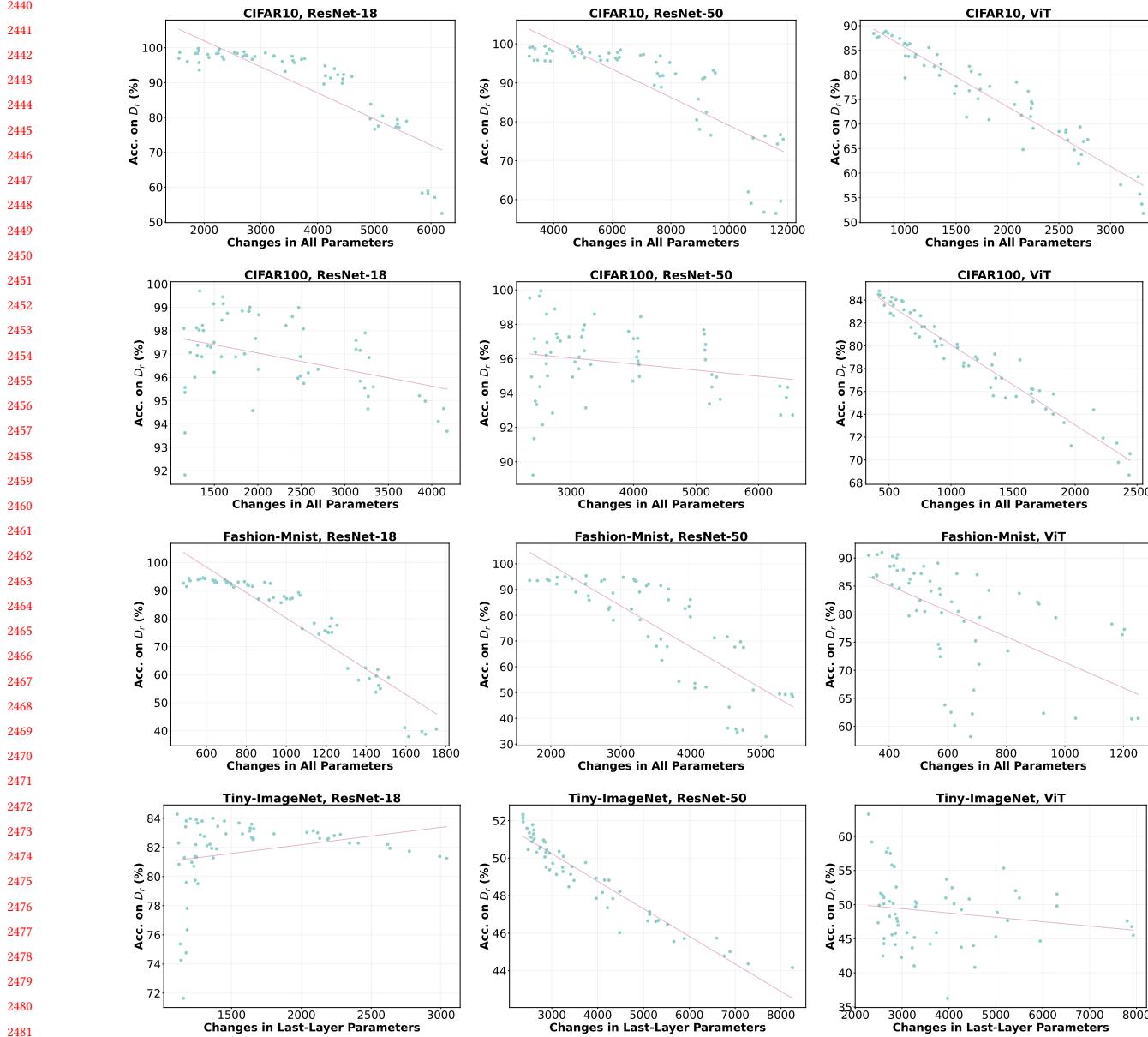


Figure 6: Dependence between the changes in all parameters and the changes in the last layer parameters by using our method. The magnitude of the parameter changes is quantified by the L_1 -norm. The result shows a positive correlation.

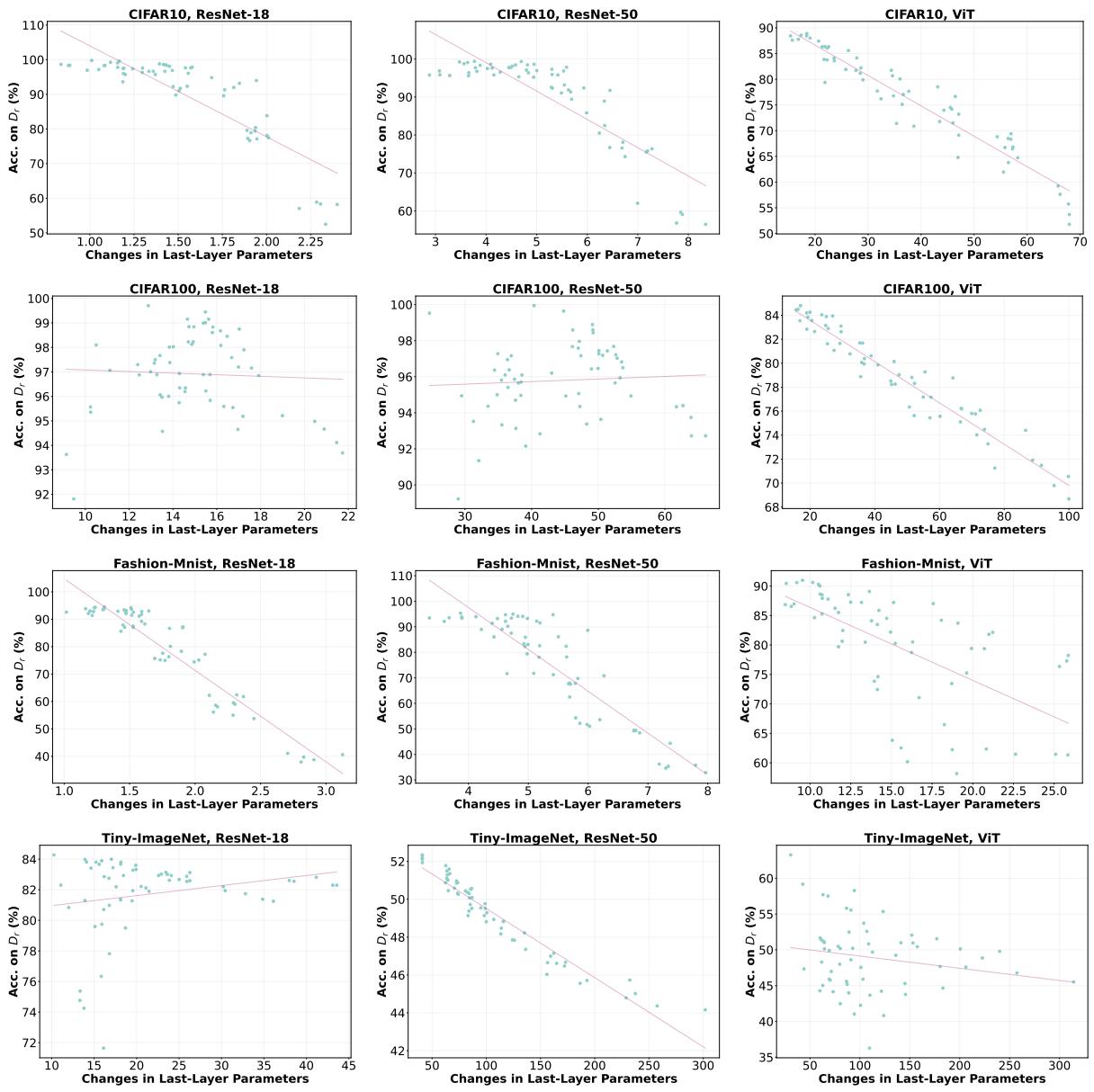
2437 **B.8 Relations between a Model's Accuracy on Remaining Examples and Changes of a Model's All
2438 Parameters.**



2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483
2484
2485
2486
2487
2488
2489
2490
2491
2492
2493
2494
2495
2496
2497
2498
2499
2500
2501
2502
2503
2504
2505
2506
2507
2508
2509
2510
2511
2512
2513
2514
2515
2516
2517
2518
2519
2520
2521
2522
2523
2524
2525
2526
2527
2528
2529
2530
2531
2532
2533
2534
2535
2536
2537
2538
2539
2540
2541
2542
2543
2544
2545
2546
2547
2548
2549
2550
2551
2552

Figure 7: Dependence between changes in all parameters and model accuracy on remaining examples by using our method. The result shows a negative correlation.

2553 **B.9 Relations between a Model's Accuracy on Remaining Examples and Changes of a Model's**
 2554 **Last-Layer Parameters.**



2559 **Figure 8: Dependence between changes in last-layer parameters and model accuracy on remaining examples by using our**
 2600 **method. The result shows a negative correlation.**