# Learning to Unlearn: Instance-wise Unlearning for Pre-trained Classifiers

Sungmin Cha[1,2]*, Sungjun Cho[2]*, Dasol Hwang[2]*,
Honglak Lee[2,4], Taesup Moon[1], and Moontae Lee[2,3]
[1] Seoul National University [2] LG AI Research
[3] University of Illinois at Chicago [4] University of Michigan
sungmin.cha@snu.ac.kr, {sungjun.cho, dasol.hwang}@lgresearch.ai,
tsmoon@snu.ac.kr, {honglak.lee, moontae.lee}@lgresearch.ai

## Abstract

*Since the recent advent of regulations for data protection (e.g., the General Data Protection Regulation), there has been increasing demand in deleting information learned from sensitive data in pre-trained models without retraining from scratch. The inherent vulnerability of neural networks towards adversarial attacks and unfairness also calls for a robust method to remove or correct information in an instance-wise fashion, while retaining the predictive performance across remaining data. To this end, we define instance-wise unlearning, of which the goal is to delete information on a set of instances from a pre-trained model, by either misclassifying each instance away from its original prediction or relabeling the instance to a different label. We also propose two methods that reduce forgetting on the remaining data: 1) utilizing adversarial examples to overcome forgetting at the representation-level and 2) leveraging weight importance metrics to pinpoint network parameters guilty of propagating unwanted information. Both methods only require the pre-trained model and data instances to forget, allowing painless application to real-life settings where the entire training set is unavailable. Through extensive experimentation on various image classification benchmarks, we show that our approach effectively preserves knowledge of remaining data while unlearning given instances in both single-task and continual unlearning scenarios.*

## 1. Introduction

Humans remember and forget: efficiently learning useful knowledge yet regulating privately sensitive information and protecting from malicious attacks. Recent advances in large-scale pre-training enable models to memorize massive information for intelligent operations [1], but there is a cost.

Vision models trained on numerous image data sometimes misclassify naturally adversarial or adversarially attacked examples with high-confidence [2]. A naïve solution is to retrain these models from scratch after refining or reweighting their training datasets [3, 4, 5]. Language models trained on indiscriminately collected data often disclose private information such as occupations, phone numbers, and family background during text generation [6]. However, such post-hoc processing is impractical due to growing volumes of data and substantial cost of large-scale training: while exercising the Right to be Forgotten [7, 8] may be straightforward to humans, it is not so straightforward in the context of machine learning. This has sparked the field of *machine unlearning*, in which the main goal is to efficiently delete information while preserving information on the remaining data.

While several machine unlearning approaches have shown promising results deleting data from traditional machine learning algorithms [9, 10, 11] as well as DNN-based classifiers [12, 13, 14, 15, 16, 17, 18], existing work are built upon assumptions far too restrictive compared to real-life scenarios. First off, many approaches assume a *class-wise unlearning* setup, where the task is to delete information from all data points that belong to a particular class or set of classes. However, data deletion requests are practically received at a per-instance basis, potentially resulting in a set of data points with a mixture of class labels [6, 19]. Another widely used assumption is that at least a subset of the original training data is available at the time of unlearning [12, 16, 17, 18]. Unfortunately, loading the original dataset may not be an option in real settings due to data expiration policies or lack of storage for large amounts of data. Lastly, many approaches consider the main objective as *undoing* the previous effect of the deleting data during training. However, previous work have shown that fulfilling this objective can still lead to information leakage, and unlearning mechanisms must *explicitly enforce misprediction* for tighter security against membership inference or model inversion attacks [20].

In light of aforementioned limitations, we propose a framework for *instance-wise unlearning* that deletes informa-

---

*Equal contribution

tion access only to the pre-trained model and the data points requested for unlearning. Instead of undoing the previous influence of deleting data, we pursue a stronger goal where all data points requested for deletion are misclassified, preventing collection of information via interpolation of nearby data points. Achieving this stronger goal for unlearning while maintaining the efficacy of the pre-trained model on the remaining data is a challenging task. Inspired by works in continual learning and adversarial attack literature [21, 22], we thus propose two regularization methods that minimize the loss in predictive performance on the remaining data. Specifically, we 1) generate adversarial examples by attacking each deleting data point with the pre-trained model and retrain on these examples to prevent representation-level forgetting and 2) use weight importance measures from unlearning instances to focus gradient updates more towards parameters responsible for the originally correct classification of such instances. Extensive experiments on CIFAR-10/-100 [23] and ImageNet-1K [24] datasets show that our proposed method effectively preserves overall predictive performance, while completely forgetting images requested for deletion. Our qualitative analyses also reveal interesting insights, including lack of any discernible pattern in misclassification that may be exploited by adversaries, preservation of the previously learned decision boundary, and forgetting of high-level features within deleted images. In summary, our **main contributions** are as follows:

- We propose *instance-wise unlearning* through intended misclassification, which can effectively unlearn data with only the pre-trained model and data to forget.

- We present two model-agnostic regularization methods that reduce forgetting on the remaining data while misclassifying data requested for deletion.

- Empirical evaluations on well-known image classification benchmarks show that our proposed method significantly boosts predictive performance after unlearning.

## 2. Related Work

**Machine unlearning.** Machine unlearning [25] is a field that makes a pre-trained model forget information learned from a specified subset of data. For this, existing studies have taken an approach that deletes the influence of unwanted data points from the model while retaining the predictive performance on the rest of the data [9, 10, 11] proposed unlearning methods for a linear/logistic regression, k-means clustering, and random forests, respectively. These methods are specifically designed for simple machine learning models, not for neural networks.

More recently, machine unlearning for deep neural networks have been studied in various settings, shown in Table 1. These methods can be categorized into two approaches:

*class-wise* and *instance-wise unlearning*. The *class-wise unlearning* is to forget all data points belonging to a certain class (*e.g.,* all images of dogs in CIFAR-10) while retaining performance on the remaining classes [12, 13, 14, 15]. On the other hand, *instance-wise unlearning* deletes information from individual data instances, possibly with mixed classes, from the pre-trained model [16, 17, 18].

Table 1: Comparison between existing unlearning methods.

| Methods | Unit | Goal | $D_r$ | $D_f$ |
|---|---|---|---|---|
| Tarun et al. [12] | class | undo | ✓ | ✗ |
| Chundawat et al. [13] | class | undo | ✗ | ✗ |
| Ye et al. [14] | class | undo | ✗ | ✓ |
| Yoon et al. [15] | class | undo | ✗ | ✓ |
| Golatkar et al. [16] | instance | undo | ✓ | ✓ |
| Kim and Woo [17] | instance | undo | ✓ | ✓ |
| Mehta et al. [18] | instance | undo | ✓ | ✓ |
| **Our methods** | **instance** | **misclassify** | ✗ | ✓ |

Previously, all existing work have set the unlearning objective of guiding the pre-training model towards the model retrained on the dataset after removing unwanted data instances (*i.e.* to *undo* their effect during training). Unfortunately, previous work have shown that this goal fails to provide reliable security against data leakage [20], due to high interpolation capability of deep neural networks. Therefore, we instead define our goal for unlearning as tuning the pre-trained model towards intentionally *misclassify* the data points requested for deletion.

Furthermore, existing methods assume different access levels to the unlearning data $D_f$ and the remaining data $D_r$. Existing solutions for *instance-wise unlearning* require access to the entire dataset (*i.e.,* $\mathcal{D}_r \cup \mathcal{D}_f$), limiting applicability in real-world scenarios where data expiration or storage issues may occur. On the other hand, our proposed methods only need to the unlearning dataset $\mathcal{D}_f$.

**Adversarial examples.** Since the vulnerability of neural networks has been revealed [26], various methods have been proposed to generate adversarial examples that can deceive neural networks [27, 28, 29, 30]. In the case of white-box attack, an adversarial example can be generated by adding a hardly visible perturbation on a given image based on the gradient information from the model, making the model classify the image to a wrong class. The injected noise of the example is hard to distinguish visually but it causes a serious misclassification of the model. Recently, [31] have experimentally demonstrated that this noise contains useful features about the attack target label for the model. Inspired by this observation, we leverage adversarial examples as a means for regularization towards preserving the previously learned decision boundary in the feature space.

**Weight importance.** Weight importance is a measure of how important each weight is when the model predicts an output for a given input data, and has been used for various
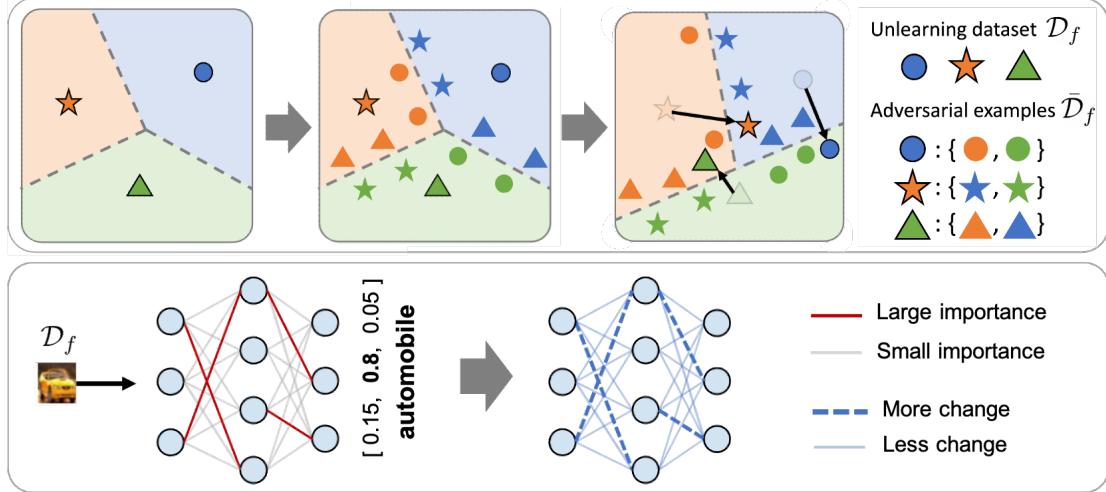
Figure 1: Illustrations of our approaches that reduce forgetting on the remaining data. (Top) Augmenting adversarial examples from unlearning data provides support for preserving the overall decision boundary. (Bottom) Weight importance measures allow us to pinpoint weights we should change to induce misclassification while maintaining other weights to mitigate forgetting.

purposes such as weight pruning [32, 33, 34, 35, 36] and regularization-based continual learning [37, 22, 38, 39, 40, 41]. Among them, *regularization-based continual learning* has actively proposed various methods for measuring the weight importance. For overcoming catastrophic forgetting of previous tasks, the weight importance is utilized as the strength of the L2 regularization between a current model's weight and the model's weight trained up to the previous task. Most methods estimate the weight-level importance based on a gradient of a given input data [37, 22].

## 3. Method

### 3.1. Preliminaries and notations

**Dataset and pre-trained model.** Let $\mathcal{D}_{train}$ be the entire training dataset used to pre-train a classification model $g_{\boldsymbol{\theta}} : \mathcal{X} \to \mathcal{Y}$. We denote $\mathcal{D}_f \subset \mathcal{D}_{train}$ as the unlearning dataset that we want to intentionally forget from the pre-trained model and $\mathcal{D}_r$ as the remaining dataset on which we wish to maintain predictive accuracy ($\mathcal{D}_r := \mathcal{D}_{train} \setminus \mathcal{D}_f$). We denote a pair of an input image $\boldsymbol{x} \in \mathcal{X}$ and its ground-truth label $\boldsymbol{y} \in \mathcal{Y}$ from $\mathcal{D}_{train}$ as $(\boldsymbol{x}, \boldsymbol{y}) \sim \mathcal{D}_{train}$, similarly $(\boldsymbol{x}_f, \boldsymbol{y}_f) \sim \mathcal{D}_f$ and $(\boldsymbol{x}_r, \boldsymbol{y}_r) \sim \mathcal{D}_r$. Also, $\mathcal{D}_{test}$ denotes the test dataset used for evaluation. Note that our approaches assumes access to only the pre-trained model $g_{\boldsymbol{\theta}}$ and the unlearning dataset $\mathcal{D}_f$ during unlearning.

**Adversarial examples.** The goal of an adversarial attack on an input $(\boldsymbol{x}, \boldsymbol{y})$ is to generate an adversarial example $\boldsymbol{x}'$ that is similar to $\boldsymbol{x}$, but leads to misclassification ($g_{\boldsymbol{\theta}}(\boldsymbol{x}') \neq \boldsymbol{y}$) when fed to the pre-trained model $g_{\boldsymbol{\theta}}$. In the case of *targeted* adversarial attack, it makes the model predict a specific class different from the true class ($g_{\boldsymbol{\theta}}(\boldsymbol{x}') = \bar{\boldsymbol{y}}$). The typical optimization form of generating adversarial examples in

targeted attack is denoted as

$$\boldsymbol{x}' = \argmin_{\boldsymbol{z}:\|\boldsymbol{z}-\boldsymbol{x}\|_p \leq \epsilon} \mathcal{L}_{\text{CE}}(g_{\boldsymbol{\theta}}(\boldsymbol{z}), \bar{\boldsymbol{y}}; \boldsymbol{\theta}) \qquad (1)$$

where $\mathcal{L}_{\text{CE}}$ stands for the cross-entropy loss. The $\|\boldsymbol{z}-\boldsymbol{x}\|_p \leq \epsilon$ condition requires that the $L_p$-norm is less than a perturbation budget $\epsilon$. The optimization above is intractable in general, and thus several papers have proposed approximations that can generate adversarial examples without directly solving it [27, 28, 30, 29]. In this paper, we make use of $L_2$-PGD targeted attacks [29] for all experiments.

**Measuring weight importance with MAS.** To measure weight importance $\Omega$, we consider MAS [22], an algorithm that estimates weight importance by finding parameters that bring a significant change in the output when perturbed slightly. It estimates the weight importance via a sum of gradients on the L2-norm of the outputs:

$$\Omega_i = \frac{1}{N} \sum_{n=1}^{N} \left| \frac{\partial \|g_{\boldsymbol{\theta}}(\boldsymbol{x}^{(n)}; \boldsymbol{\theta})\|_2^2}{\partial \theta_i} \right| \qquad (2)$$

where $i$ stands for the index of network parameter weights and $\boldsymbol{x}^{(n)}$ denotes $n$-th input image from a total of $N$ numbers of images. Note that each $\Omega_i$ can be interpreted as a measure of influence or importance of $\theta_i$ in producing the output of given $N$ input images.

### 3.2. Our proposed framework

**Definition of instance-wise unlearning.** Let $\hat{g}_{\boldsymbol{\theta}}$ denote the model after unlearning. We consider two types of misclassification for instance-wise unlearning: (**i**) *misclassifying* all data points in $\mathcal{D}_f$, (*i.e.*, $\hat{g}_{\boldsymbol{\theta}}(\boldsymbol{x}_f) \neq \boldsymbol{y}_f$). (**ii**) *relabeling (or*

**Algorithm 1** Generate adversarial examples
---
**Require:** Forgetting data $\mathcal{D}_f$, Model $g_{\boldsymbol{\theta}}$
**Ensure:** Adversarial examples $\bar{\mathcal{D}}_r$
1: $\bar{\mathcal{D}}_r \leftarrow \emptyset$
2: **for** $i$ **in range** $N_f$ **do**
3: $\quad (\boldsymbol{x}^{(i)}, \boldsymbol{y}^{(i)}) \sim \mathcal{D}_f$
4: $\quad$ Randomly sample $\bar{\boldsymbol{y}}^{(i)} \neq \boldsymbol{y}^{(i)}$
5: $\quad$ **for** $j$ **in range** $N_{adv}$ **do**
6: $\quad\quad \boldsymbol{x}_f'^{(j)} \leftarrow L_2\text{-PGD}(\boldsymbol{x}^{(i)}, \bar{\boldsymbol{y}}^{(i)})$ (Eq. 1)
7: $\quad\quad \bar{\mathcal{D}}_r \leftarrow \bar{\mathcal{D}}_r \cup \{(\boldsymbol{x}_f'^{(j)}, \bar{\boldsymbol{y}}^{(j)})\}$
8: $\quad$ **end for**
9: **end for**
10: **return** $\bar{\mathcal{D}}_r$

**Algorithm 2** Measure weight importance
---
**Require:** Forgetting data $\mathcal{D}_f$, Model $g_{\boldsymbol{\theta}}$
**Ensure:** Weight importance $\bar{\Omega}$
1: $\bar{\Omega} \leftarrow \{0\}$
2: $\Omega \leftarrow$ weight importances$(\mathcal{D}_f, g_{\boldsymbol{\theta}})$ (Eq. 2)
3: **for** $l$ **in range** $L$ **do**
4: $\quad$ Get importance of $l$-th layer $\Omega^l \leftarrow \Omega$
5: $\quad$ Normalize $\Omega^l \leftarrow \dfrac{\Omega^l - Min(\Omega^l)}{Max(\Omega^l) - Min(\Omega^l)}$
6: $\quad$ Update $\bar{\Omega}^l \leftarrow \{1 - \Omega^l\}$
7: **end for**
8: **return** $\bar{\Omega}$

*correcting)* the predictions of $\mathcal{D}_f$ (*i.e.*, $\hat{g}_{\boldsymbol{\theta}}(\boldsymbol{x}_f) = \boldsymbol{y}_f^*$) where $\boldsymbol{y}_f^* \neq \boldsymbol{y}_f$ is chosen individually for each input $\boldsymbol{x}_f$. Let $\mathcal{L}_{\text{UL}}$ denote a loss function used for unlearning on a classification model. The above two goals can be realized with the following loss functions:

$$\mathcal{L}_{\text{UL}}^{\text{MS}}(D_f; \boldsymbol{\theta}) = -\mathcal{L}_{\text{CE}}(g_{\boldsymbol{\theta}}(\boldsymbol{x}_f), \boldsymbol{y}_f; \boldsymbol{\theta}) \qquad (3)$$

$$\mathcal{L}_{\text{UL}}^{\text{Cor}}(D_f; \boldsymbol{\theta}) = \mathcal{L}_{\text{CE}}(g_{\boldsymbol{\theta}}(\boldsymbol{x}_f), \boldsymbol{y}_f^*; \boldsymbol{\theta}) \qquad (4)$$

When unlearning solely based on the two loss functions above, the model is likely to suffer from significant forgetting on $\mathcal{D}_r$. Therefore, a crucial objective shared across both unlearning goals is to overcome forgetting of previously learned knowledge, and maintain as much classification accuracy as possible on $\mathcal{D}_r$.

When both $\mathcal{D}_f$ and $\mathcal{D}_r$ are available, we can easily obtain an *oracle* model that satisfies the objective by re-training the model with the following loss function: $\mathcal{L}_{\text{oracle}}(\mathcal{D}_f, \mathcal{D}_r; \boldsymbol{\theta}) = \mathcal{L}_{\text{UL}}(\mathcal{D}_f; \boldsymbol{\theta}) + \mathcal{L}_{\text{CE}}(\mathcal{D}_r; \boldsymbol{\theta})$. However in real-settings, access to $D_r$ may not be an option due to high cost in data storage. To tackle this limitation, we define a regularization-based unlearning that achieves the goal above without use of $\mathcal{D}_r$:

$$\mathcal{L}_{\text{RegUL}}(D_f; \boldsymbol{\theta}) = \mathcal{L}_{\text{UL}}(D_f; \boldsymbol{\theta}) + \mathcal{R}(D_f, g_{\boldsymbol{\theta}}) \qquad (5)$$

Here, $\mathcal{R}(\cdot)$ is the regularization term used to overcome forgetting of knowledge on the remaining data $\mathcal{D}_r$. In the following subsections, we introduce two novel regularization methods designed to overcome representation- and weight-level forgetting during the unlearning process.

**Regularization using adversarial examples.** The motivation of using adversarial examples stems from the work of [31], which showed that perturbations added to $\boldsymbol{x}$ to generate an adversarial example $\boldsymbol{x}'$ contain class-specific features of the attack target label $\bar{\boldsymbol{y}} \neq \boldsymbol{y}$. Based on this finding, we utilize generated adversarial examples as part of regularization $\mathcal{R}(\cdot)$ to preserve class-specific knowledge previously learned by the model, overcoming forgetting during

unlearning at the *representation-level*. Let $D_f$ be a set of $N_f$ images: $\{(\boldsymbol{x}_f^{(i)}, \boldsymbol{y}_f^{(i)})\}_{i=1}^{N_f}$. Prior to the unlearning process, we generate adversarial examples $\boldsymbol{x}_f'$ using the targeted PGD attack with a randomly selected attack target label $\bar{\boldsymbol{y}} \neq \boldsymbol{y}_f$. We generate $N_{\text{adv}}$ adversarial examples per input $\boldsymbol{x}_f$, resulting in a set of $\bar{N}_f = N_f \times N_{\text{adv}}$ examples, denoted as $\bar{\mathcal{D}}_f = \{(\boldsymbol{x}_f'^{(i)}, \bar{\boldsymbol{y}}_f^{(i)})\}_{i=1}^{\bar{N}_f}$. During unlearning, we add $\mathcal{L}_{\text{CE}}(\bar{\mathcal{D}}_f; \boldsymbol{\theta})$ as a regularization term with adversarial examples:

$$\begin{aligned} \mathcal{L}_{\text{UL}}^{\text{Adv}}(D_f; \boldsymbol{\theta}) &= \mathcal{L}_{\text{UL}}(D_f; \boldsymbol{\theta}) + \mathcal{R}_{\text{Adv}}(D_f, g_{\boldsymbol{\theta}}) \\ &= \mathcal{L}_{\text{UL}}(D_f; \boldsymbol{\theta}) + \mathcal{L}_{\text{CE}}(\bar{\mathcal{D}}_f; \boldsymbol{\theta}) \end{aligned} \qquad (6)$$

An intuitive illustration of this approach in the representation-level is shown in Figure 1. In essence, the generated adversarial examples in $\bar{\mathcal{D}}_f$ mimic the remaining dataset $\mathcal{D}_r$, providing information of the pre-trained decision boundary within the representation space. As a result, by adding $\mathcal{L}_{\text{CE}}(\bar{\mathcal{D}}_f; \boldsymbol{\theta})$ as a regularizer to the unlearning process, the model learns a new decision boundary that minimizes $\mathcal{L}_{\text{UL}}$ (in Eq. 3 and 4) while simultaneously attempting to keep the decision boundary of the original model. The pseudocode for generating adversarial examples is in Algorithm 1.

**Regularization with weight importance.** We also propose a regularization using weight importance to overcome forgetting at the weight-level. As depicted in Figure 1, our approach is to maintain the weights that were less important for predictions on $\mathcal{D}_f$ as much as possible, while allowing changes in weights that are considered important for correctly predicting $\mathcal{D}_f$. In other words, we penalize changes of weights that were less important when classifying $\mathcal{D}_f$, thereby preventing weight-level forgetting.

Specifically, we calculate the weight importance with MAS given $g_{\boldsymbol{\theta}}$ and $\mathcal{D}_f$ before unlearning, and normalize the importance measurements $\Omega^l$ within each $l$-th layer to lie within $[0, 1]$. Note that this normalized importance $\Omega^l$ assigns large values to weights important for classifying

$\mathcal{D}_f$. Therefore, we add a regularization term that sums all parameter change in the $l$-th layer weighted by $\bar{\Omega}^l = 1 - \Omega^l$, such that more important weights are updated more. The objective including weight importance regularization with regularization via adversarial examples can be written as:

$$\begin{aligned}\mathcal{L}_{\text{UL}}^{\text{Adv+Imp}}(D_f; \boldsymbol{\theta}) &= \mathcal{L}_{\text{UL}}^{\text{Adv}}(D_f; \boldsymbol{\theta}) + \mathcal{R}_{\text{Imp}}(D_f, g_{\boldsymbol{\theta}}) \\ &= \mathcal{L}_{\text{UL}}^{\text{Adv}}(D_f; \boldsymbol{\theta}) + \sum_i \bar{\Omega}_i(\theta_i - \tilde{\theta}_i)^2 \end{aligned} \quad (7)$$

where $i$ is the index of each weight and $\tilde{\theta}$ is the initial weight of the pre-trained classifier before unlearning. The pseudocode of measuring weight importance is shown in Algorithm 2. Throughout various experiments, we observe that applying the regularization using adversarial examples is already effective to overcome the forgetting for knowledge of $\mathcal{D}_r$, and the additional regularization with weight importance further enhances performance even further, especially in more challenging scenarios such as continual unlearning. The pseudocode of our overall unlearning pipeline can be found in the supplementary material.

# 4. Experiments

In this section, we evaluate our method in various image classification benchmarks. We first describe our experimental setup, including datasets, baselines, and experimental details. We then show that our methods effectively preserve the knowledge of remaining data while unlearning instances that should be forgotten in both single-task and continual unlearning scenarios. Lastly, we offer qualitative analyses on three different aspects: prediction patterns, decision boundaries, and layer-wise representations in unlearning.

## 4.1. Setup

**Datasets and baselines.** We evaluate our unlearning methods on three different image classification datasets: CIFAR-10, CIFAR-100 [23], and ImageNet-1K [24]. We use the ResNet-18 and ResNet-50 [42] as the base model for CIFAR-10 and CIFAR-100/ImageNet-1K, respectively. The experimental results from various base models are also available in the supplementary material. The baselines are as follows: BEFORE corresponds to the pre-trained model before unlearning; ORACLE denotes the ideal model that is retrained with positive gradients from $\mathcal{D}_r$ and negative gradients from $\mathcal{D}_f$ until maximum performance is achieved on $\mathcal{D}_r$ (i.e., $\mathcal{L}_{\text{oracle}}$). In case the accuracy of $\mathcal{D}_f$ does not reach 0 after a certain number of epochs, we report results based on the checkpoint that attained the lowest performance on $\mathcal{D}_f$; Repeated adversarial weight perturbation (RAWP) is another baseline derived from AWP [43] that induces misclassification by adding weight perturbations generated from $\mathcal{D}_f$ to the model parameters. More details on RAWP can be found in the supplementary materials. NEGGRAD [16] fine-tunes

the model on $\mathcal{D}_f$ using negative gradients (i.e., $\mathcal{L}_{\text{UL}}^{\text{MS}}$) without any regularization; CORRECT is analogous to NEGGRAD, but uses relabeling to certain label instead (i.e., $\mathcal{L}_{\text{UL}}^{\text{Cor}}$); ADV denotes our proposed method using adversarial examples (i.e., $\mathcal{L}_{\text{UL}}^{\text{Adv}}$); Lastly, ADV+IMP uses adversarial examples as well as weight importance for regularization (i.e., $\mathcal{L}_{\text{UL}}^{\text{Adv+Imp}}$).

**Experimental details.** For each dataset, we randomly pick $k \in \{16, 64, 128, 256\}$ images from the entire training dataset as unlearning data $\mathcal{D}_f$ and consider the remaining data as $\mathcal{D}_r$. For unlearning, we use a SGD optimizer with a learning rate of 1e-3, weight decay of 1e-5, and momentum of 0.9 across all experiments. We early-stop training when the model attains 0% or 100% accuracy on the unlearning data $\mathcal{D}_f$, in case of *misclassifying* and *relabeling*, respectively. For generating adversarial examples from $\mathcal{D}_f$, we use $L_2$-PGD targeted attack [29] with a learning rate of 1e-1, attack iterations of 100 and $\epsilon = 0.4$. It generates 20 adversarial examples for CIFAR-10 and 200 examples for CIFAR-100 and ImageNet-1K. For the weight importance regularization, we set regularization strength $\lambda = 1$ in Eq. 5.

## 4.2. Main Results

**Results on various datasets.** Table 2 shows evaluation results before and after unlearning $k$ instances from ResNet-18 and ResNet-50 models pre-trained on CIFAR-10 and CIFAR-100/ImageNet-1K, respectively. With respect to accuracies on $\mathcal{D}_f$, we find that all methods can almost forget up to $k = 256$ instances with consistently zero post-unlearning accuracies, except for some cases (e.g., ORACLE, NEGGRAD and RAWP in CIFAR-10 when $k = 256$). For most cases, ORACLE shows not only superior accuracy in the remaining data (i.e. $\mathcal{D}_r$ and $\mathcal{D}_{test}$) but also almost zero accuracy in $D_f$. Nevertheless, we confirm that, the size of $k$ increases, its performance decreases, demonstrating the difficulty of achieving the proposed goal of unlearning. NEGGRAD results in a significant loss of accuracy on the remaining data for all datasets. RAWP achieves competitive accuracy in the remaining data than NEGGRAD in CIFAR-10. However, it shows significantly worse performance when the dataset become complex (e.g., CIFAR-100 and ImageNet) or $k$ become more larger than $k = 16$, demonstrating simply perturbing the parameters cannot be a direct solution. Meanwhile, adding regularization with adversarial examples (ADV) boosts the accuracy by more than 40% depending on the number of instances to forget. Incorporating weight importances from MAS (ADV+IMP) provides further improvement, even supassing the performance of ORACLE in some cases (e.g., $k = 256$ in CIFAR-10 and ImageNet).

Table 5 shows results analogous to Table 2, but with the goal of relabeling data points in $\mathcal{D}_f$ to arbitrarily chosen labels rather than misclassifying. We find that a similar trend, first, ORACLE achieves superior accuracy in most cases but we observe that it suffers from the trade-off between $D_f$ and

Table 2: Evaluation results before and after unlearning $k$ instances from ResNet-18 (for CIFAR-10) and ResNet-50 (for CIFAR-100 and ImageNet-1K) pretrained on respective image classification datasets. While using negative gradients only loses significant information on $\mathcal{D}_r$, our proposed methods ADV and ADV+IMP retain predictive performance on $\mathcal{D}_r$ as well as $\mathcal{D}_{test}$, while completely forgetting instances in $\mathcal{D}_f$.

| | | CIFAR-10 | | | | CIFAR-100 | | | | ImageNet-1K | | | |
| | | $k=16$ | $k=64$ | $k=128$ | $k=256$ | $k=16$ | $k=64$ | $k=128$ | $k=256$ | $k=16$ | $k=64$ | $k=128$ | $k=256$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $D_f$ ($\downarrow$) | BEFORE | 100.0 | 99.38 | 99.53 | 99.38 | 100.0 | 100.0 | 100.0 | 100.0 | 87.50 | 84.69 | 86.09 | 86.02 |
| | ORACLE | 0.0 | 0.0 | 0.0 | 0.63 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.26 | 5.73 |
| | NEGGRAD | 0.0 | 0.0 | 0.0 | 3.83 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| | RAWP | 0.0 | 0.0 | 4.06 | 6.17 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| | ADV | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| | ADV+IMP | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 2.73 |
| $D_r$ ($\uparrow$) | BEFORE | 99.60 | 99.60 | 99.60 | 99.60 | 99.98 | 99.98 | 99.98 | 99.98 | 87.42 | 87.42 | 87.42 | 87.42 |
| | ORACLE | 98.74 | 99.72 | 98.97 | 39.90 | 99.96 | 96.17 | 96.74 | 96.43 | 73.36 | 77.81 | 67.7 | 61.61 |
| | NEGGRAD | 15.79 | 9.22 | 7.11 | 6.34 | 66.97 | 26.20 | 11.64 | 1.70 | 65.02 | 35.66 | 19.27 | 0.48 |
| | RAWP | 67.15 | 37.32 | 13.28 | 9.12 | 72.44 | 17.39 | 2.82 | 1.12 | 20.33 | 2.72 | 0.29 | 0.10 |
| | ADV | 69.70 | 66.97 | 53.49 | 39.33 | 89.18 | 81.07 | 76.28 | 65.67 | 80.92 | 73.60 | 66.31 | 43.70 |
| | ADV+IMP | **85.75** | **72.77** | **54.51** | **48.95** | **89.91** | **89.48** | **82.86** | **67.60** | **81.43** | **79.03** | **74.14** | **65.32** |
| $D_{test}$ ($\uparrow$) | BEFORE | 92.59 | 92.59 | 92.59 | 92.59 | 77.10 | 77.10 | 77.10 | 77.10 | 76.01 | 76.01 | 76.01 | 76.01 |
| | ORACLE | 90.21 | 91.01 | 89.44 | 38.59 | 64.41 | 67.06 | 66.88 | 65.31 | 63.91 | 68.08 | 59.09 | 54.03 |
| | NEGGRAD | 15.87 | 9.28 | 7.11 | 6.47 | 48.07 | 21.11 | 10.19 | 1.71 | 56.70 | 31.10 | 17.08 | 0.46 |
| | RAWP | 63.10 | 35.47 | 12.91 | 9.12 | 50.99 | 14.41 | 2.79 | 1.12 | 17.53 | 2.41 | 0.28 | 0.10 |
| | ADV | 65.14 | 62.23 | 49.47 | 36.69 | 63.17 | 57.43 | 53.89 | 46.45 | 70.36 | 65.51 | 57.34 | 38.25 |
| | ADV+IMP | **79.65** | **67.08** | **50.82** | **45.44** | **63.69** | **62.83** | **58.44** | **48.71** | **70.77** | **68.43** | **64.05** | **56.52** |

Table 3: Correcting adversarial images from ImageNet-A. ADV achieves the least forgetting, while ADV+IMP fails to correct large number of predictions due to strong regularization.

| | | ImageNet-A | | | |
| | | $k=16$ | $k=32$ | $k=64$ | $k=128$ |
|---|---|---|---|---|---|
| $D_f$ ($\uparrow$) | BEFORE | 0.0 | 0.0 | 0.0 | 0.0 |
| | CORRECT | 100.0 | 100.0 | 100.0 | 100.0 |
| | ADV | 100.0 | 100.0 | 95.31 | 83.44 |
| | ADV+IMP | 100.0 | 100.0 | 10.94 | 9.38 |
| $D_r$ ($\uparrow$) | BEFORE | 87.46 | 87.46 | 87.46 | 87.46 |
| | CORRECT | **84.41** | 83.29 | 80.79 | 77.38 |
| | ADV | 81.75 | **83.80** | **83.74** | **83.44** |
| | ADV+IMP | 81.82 | 83.73 | 83.53 | 82.86 |
| $D_{test}$ ($\uparrow$) | BEFORE | 76.15 | 76.15 | 76.15 | 76.15 |
| | CORRECT | **73.21** | 72.04 | 69.91 | 66.73 |
| | ADV | 70.89 | **72.58** | **72.68** | **72.36** |
| | ADV+IMP | 70.98 | 72.51 | 72.39 | 71.68 |

Table 4: Unlearning instances continually by increments of $k_{CL} = 8$ images per step. Our methods outperform NEGGRAD in the continual unlearning scenario as well.

| | | CIFAR-100 ($k_{CL}=8$) | | | |
| | | $k=8$ | $k=16$ | $k=64$ | $k=128$ |
|---|---|---|---|---|---|
| $D_f$ ($\downarrow$) | BEFORE | 100.0 | 100.0 | 100.0 | 100.0 |
| | NEGGRAD | 0.0 | 0.0 | 0.0 | 0.52 |
| | ADV | 0.0 | 0.0 | 1.04 | 0.0 |
| | ADV+IMP | 0.0 | 0.0 | 0.0 | 1.04 |
| $D_r$ ($\uparrow$) | BEFORE | 99.98 | 99.98 | 99.98 | 99.98 |
| | NEGGRAD | 80.58 | 31.85 | 6.60 | 1.89 |
| | ADV | 80.33 | 70.54 | 59.67 | 38.16 |
| | ADV+IMP | **81.46** | **72.78** | **62.30** | **47.14** |
| $D_{test}$ ($\uparrow$) | BEFORE | 77.10 | 77.10 | 77.10 | 77.10 |
| | NEGGRAD | 58.20 | 24.48 | 5.73 | 1.22 |
| | ADV | 57.56 | 50.43 | 43.48 | 30.10 |
| | ADV+IMP | **58.33** | **51.97** | **45.09** | **36.17** |

the remaining datasets when the size of $k$ become large (*e.g.*, $k = 256$). Second, RAWP shows much better performance than CORRECT in CIFAR-10, however, not for other datasets. Third, ADV attains significantly less forgetting in $\mathcal{D}_r$ and $\mathcal{D}_{test}$ compared to CORRECT, while succesfully relabeling all points in most cases. On the other hand, while ADV+IMP show even less forgetting, it loses accuracy in relabeling $\mathcal{D}_f$, showing that regularization via weight importance focuses too much on retaining previous knowledge rather than adapting to corrections provided in $\mathcal{D}_f$. An intuitive explanation on why this occurs particularly in relabeling is that while misclassifying can be done easily by driving the input to its closest decision boundary, relabeling can be difficult if the new class is far from the original class in the representation space. The difficulty rises even more when the size of $\mathcal{D}_f$ is large, in which case more parameters in the network are discouraged from being updated during unlearning.

**Correcting natural adversarial examples.** Leveraging the ImageNet-A [2] dataset consisting of natural images that are misclassified with high-confidence by strong classifiers, we test whether our method can make corrections on these adversarial examples, while preserving knowledge from the original training data. For this experiment, we consider $\mathcal{D}_f$ to consist $k$ adversarial images from ImageNet-A, and adjust a ResNet-50 model pre-trained on ImageNet-1K to correctly classify $\mathcal{D}_f$ via our unlearning framework. In Table 3, we find that correcting predictions of a small number of images (e.g. $k = 16$), finetuning the model naïvely with cross-

Table 5: Results analogous to Table 2, but with unlearning via relabeling each image in $\mathcal{D}_f$ to an arbitrarily chosen class. We see a similar trend where CORRECT loses significant information on $\mathcal{D}_r$, while our proposed methods retain predictive performance on $\mathcal{D}_r$ as well as $\mathcal{D}_{test}$.

| | | CIFAR-10 | | | | CIFAR-100 | | | | ImageNet-1K | | | |
| | | $k=16$ | $k=64$ | $k=128$ | $k=256$ | $k=16$ | $k=64$ | $k=128$ | $k=256$ | $k=16$ | $k=64$ | $k=128$ | $k=256$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | BEFORE | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| | ORACLE | 100.0 | 100.0 | 100.0 | 77.34 | 100.0 | 100.0 | 100.0 | 79.92 | 100.0 | 100.0 | 99.48 | 95.83 |
| $D_f$ (↑) | CORRECT | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 99.84 | 99.84 | 100.0 | 100.0 | 100.0 | 100.0 |
| | RAWP | 100.0 | 100.0 | 100.0 | 96.41 | 100.0 | 99.38 | 99.06 | 93.75 | 100.0 | 98.96 | 56.25 | 39.06 |
| | ADV | 100.0 | 99.38 | 98.28 | 96.17 | 100.0 | 100.0 | 98.28 | 92.42 | 100.0 | 89.06 | 71.88 | 58.59 |
| | ADV+IMP | 100.0 | 53.75 | 50.16 | 59.06 | 86.25 | 20.63 | 15.16 | 9.69 | 100.0 | 9.38 | 3.91 | 2.86 |
| | BEFORE | 99.60 | 99.60 | 99.60 | 99.60 | 99.98 | 99.98 | 99.98 | 99.98 | 87.48 | 87.42 | 87.42 | 87.42 |
| | ORACLE | 99.73 | 99.94 | 99.90 | 95.81 | 99.90 | 99.97 | 99.79 | 99.50 | 83.93 | 83.18 | 84.97 | 83.61 |
| $D_r$ (↑) | CORRECT | 11.75 | 12.33 | 9.71 | 8.97 | 74.84 | 31.79 | 18.64 | 6.93 | 82.47 | 77.21 | 68.19 | 44.47 |
| | RAWP | 79.94 | 63.95 | 52.12 | 21.14 | 59.36 | 40.93 | 9.20 | 6.15 | 48.66 | 31.55 | 0.46 | 0.24 |
| | ADV | 85.53 | 83.36 | 81.06 | 72.41 | **92.94** | 94.64 | 96.32 | **94.28** | 83.39 | **84.79** | **84.52** | **84.11** |
| | ADV+IMP | **91.15** | **94.76** | **90.57** | **81.82** | 90.73 | **96.68** | **96.44** | 92.40 | **83.58** | 83.31 | 80.43 | 79.29 |
| | BEFORE | 92.59 | 92.59 | 92.59 | 92.59 | 77.10 | 77.10 | 77.10 | 77.10 | 76.01 | 76.01 | 76.01 | 76.01 |
| | ORACLE | 91.65 | 91.99 | 91.57 | 87.51 | 74.05 | 74.93 | 74.15 | 73.81 | 72.76 | 72.04 | 73.83 | 72.79 |
| $D_{test}$ (↑) | CORRECT | 11.79 | 12.16 | 9.80 | 8.87 | 53.11 | 24.37 | 14.64 | 5.95 | 71.57 | 66.61 | 58.37 | 38.29 |
| | RAWP | 73.58 | 59.48 | 48.02 | 20.25 | 41.93 | 29.51 | 7.88 | 5.71 | 42.40 | 27.52 | 0.48 | 0.23 |
| | ADV | 79.15 | 76.95 | 74.61 | 66.80 | **65.62** | 66.79 | 68.56 | **66.75** | 72.34 | **73.51** | **73.38** | **73.00** |
| | ADV+IMP | **84.24** | **86.92** | **82.82** | **74.68** | 64.28 | **69.15** | **68.60** | 64.81 | **72.49** | 71.92 | 69.20 | 68.46 |

entropy only attains the best accuracy in both $\mathcal{D}_r$ and $\mathcal{D}_{test}$. When correcting larger number of images, however, the absence of regularization terms results in larger forgetting in $\mathcal{D}_r$ compared to ADV and ADV+IMP, with a performance gap that consistently increases with the number of adversarial images. Another takeaway is that regularization via weight importance does not help in this scenario, even showing a significant drop in $\mathcal{D}_f$ accuracy when a large number of adversarial images are introduced. This implies that using weight importances imposes too strong a regularzation that correcting predictions for $\mathcal{D}_f$ itself becomes non-trivial. We conjecture that the aggregation of important parameters for predictions in $\mathcal{D}_f$ cover a large proportion of the network with large $k$, and that careful search for the Pareto optimal between accuracies on $\mathcal{D}_f$ and on $\mathcal{D}_r$ is required.

**Continual unlearning.** In real-world scenarios, it is likely that data removal requests come as a stream, rather than all at once. Ultimately, despite continual unlearning requests, we need the unlearning method that can delete the requested data while maintaining performance for the rest data. Thus, we consider the setting of deleting $k = \{8, 16, 64, 128\}$ data by repeating the procedure of continually unlearning $\mathcal{D}_f$ in small fragments of size $k_{CL} = 8$. Table 4 shows the results of continual unlearning in the model trained with ResNet-50 on CIFAR-100. We observe that NEGGRAD suffers from large forgetting as the iteration of unlearning procedure increases. On the other hand, our proposed method shows significantly less forgetting while effectively deleting for $\mathcal{D}_f$ even after multiple iterations of unlearning.

### 4.3. Qualitative Analysis

Through further analysis, we gather insight on the following questions: **Q1.** Is there any particular pattern in how the model unlearns a set of instances (*i.e.,* does the model use
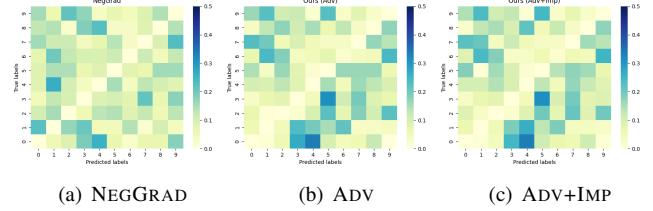


(a) NEGGRAD   (b) ADV   (c) ADV+IMP

Figure 2: Confusion matrices showing average pairwise frequencies of pre- (Y-axis) and post-unlearning (X-axis) prediction labels from $\mathcal{D}_f$. A hue closer to blue indicates higher frequency. Our unlearning framework does not produce any discernible correlation in misclassification.

any particular label as a retainer for deleted data)? **Q2.** How does the model isolate out instances in $\mathcal{D}_f$ from its previous decision boundary? **Q3.** How do layer-wise representations of data points in $\mathcal{D}_f$ and $\mathcal{D}_r$ change before and after unlearning? For interpretable visualizations, we perform the following analysis on a ResNet-18 model pre-trained on CIFAR-10.

**A1. Our method shows no pattern in misclassification.** We first check whether the unlearned model classifies all instances in $\mathcal{D}_f$ to a particular set of labels. The model exhibiting no correlation between true labels and new misclassified labels is crucial with respect to data privacy, as it indicates that the unlearning process avoids the so-called *Streisand effect* where data instances being forgotten unintentionally becomes more noticeable [16]. Figure 2 shows the confusion matrices of (pre-unlearning label, post-unlearning label) pairs from $\mathcal{D}_f$ for $k = 512$. We find no distinguishable pattern when unlearning with our methods as well as Neg-Grad, which shows that no specific label is used as a retainer, which adds another layer of security against adversaries in
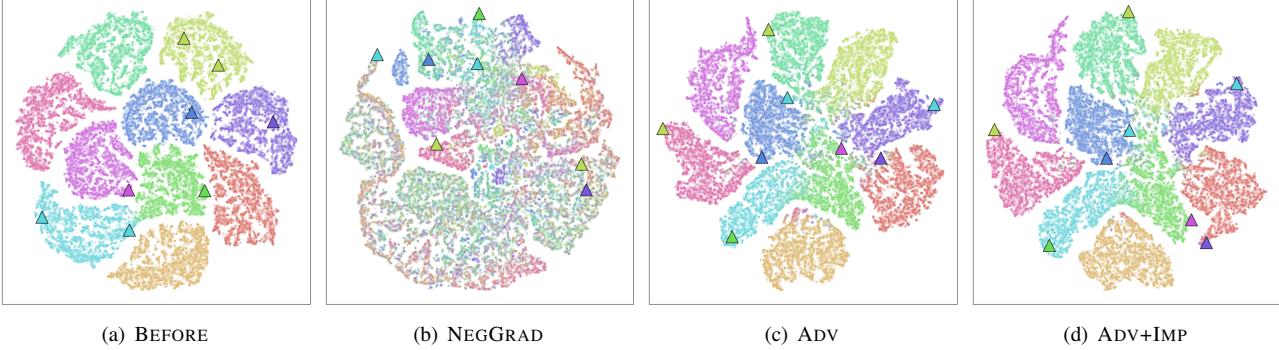
(a) BEFORE      (b) NEGGRAD      (c) ADV      (d) ADV+IMP

Figure 3: t-SNE plots of CIFAR-10 datapoints in $\mathcal{D}_f$ (triangles) and $\mathcal{D}_r$ (dots) before and after unlearning. Colors indicate true labels for all plots. Regularization with adversarial examples and weight importance effectively preserves the decision boundary while migrating instances in $\mathcal{D}_f$ towards the class boundary to induce misclassification.
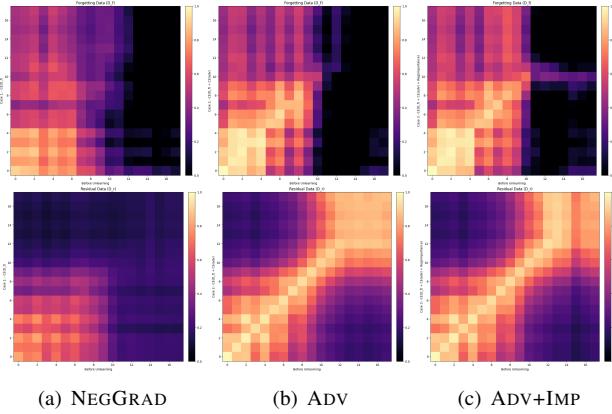


(a) NEGGRAD      (b) ADV      (c) ADV+IMP

Figure 4: Layer-wise CKA correlations on $\mathcal{D}_f$ (top row) and $\mathcal{D}_r$ (bottom row) between representations before (X-axis) and after (Y-axis) unlearning. Brighter color indicates higher CKA correlation. NEGGRAD results in large forgetting of high-level features in not only $\mathcal{D}_f$, but also $\mathcal{D}_r$. Our approaches, on the other hand, selectively forget high-level features only in $\mathcal{D}_f$.

search of unlearned data points.

**A2. Our method effectively preserves the decision boundary.** We check whether the adversarial examples generated from forgetting data help in preserving the decision boundary in the feature space. Figure 3 shows t-SNE [44] visualizations of final-layer activations from examples in $\mathcal{D}_r$ and $\mathcal{D}_f$ before and after unlearning. We find that unlearning through only negative gradient significantly distorts the previous decision boundary, leading to poor predictive performance across $\mathcal{D}_r$. However, when we incorporate adversarial samples from instances in $\mathcal{D}_f$, the decision boundary is well-preserved with unlearned examples being inferred as boundary cases in-between multiple classes. Even for examples that lie far from the decision boundary before unlearning, our method successfully relocates the corresponding representations towards the decision boundary, while keeping each

class cluster intact.

**A3. Our method unlearns data by forgetting high-level features.** Lastly, we compare the representations at each layer of the model before and after unlearning to identify where the intended forgetting occurs. For this analysis, we leverage CKA [45] which measures correlations between representations given two distinct models. Figure 4 shows the CKA correlation heatmaps between the original pre-trained model and the same model after unlearning. Results show that for examples in $\mathcal{D}_f$, representations are no longer aligned starting from the 10-th layer while the representations before that layer still resemble those from the original model. This indicates that the model forgets examples by forgetting high-level features, while similarly recognizing low-level features in images as the original model. This insight is consistent with previous observations in the continual learning literature that more forgettable examples exhibit peculiarities in high-level features [46].

## 5. Concluding Remarks

We propose an instance-wise unlearning framework that deletes information from a pre-trained model given a set of data instances with mixed labels. Rather than undoing the influence of given instances during the pre-training, we aim for a stronger form of unlearning via intended two types of misclassification, misclassifying and relabeling. We develop two regularization techniques that reduce forgetting the remaining data, one utilizing adversarial examples of deleting instances and another leveraging weight importance to focus updates on parameters responsible for propagating information we wish to forget. Both approaches are agnostic to the choice of architecture and require access only to the pre-trained model and instances requested for deletion. Experiments on various image classification datasets showed that our methods effectively mitigate forgetting on remaining data, while completely misclassifying deletion data.

As future work, we hope to further explore the efficacy of

our approach on real-world datasets that require unlearning for protection against privacy violations or unfair predictions due to gender or racial biases. Unlearning information of certain data points from a pre-trained generative model such as diffusion-based models would also be an exciting future direction.

## A. Pseudo code of overall unlearning process

---

**Algorithm 3** The pseudo-code of overall unlearning process the case of using $\mathcal{L}_{\text{UL}}^{\text{MS}}$.

---
1: UNLEARNACC = 100
2: MAXEP = 100
3: EP = 0
4: $\bar{\mathcal{D}}_r \leftarrow$ Generate adversarial examples with Algorithm 1
5: $\bar{\Omega} \leftarrow$ Measure weight importance with Algorithm 2
6: $\tilde{\boldsymbol{\theta}} \leftarrow \boldsymbol{\theta}$ // Set the pretrained model for unlearning
7: **while** UNLEARNACC $\neq 0$ **do**
8:    Minimize Eqn (6) or (7)
9:    UNLEARNACC = GetAccuracy($\mathcal{D}_f, g_{\tilde{\boldsymbol{\theta}}}$)
10:   **if** EP > MAXEP **then**
11:      break
12:      EP += 1
13:   **end if**
14: **end while**
15: **return** $\tilde{\boldsymbol{\theta}}$

---

**Algorithm 4** The pseudo-code of overall unlearning process the case of using $\mathcal{L}_{\text{UL}}^{\text{Cor}}$.

---
1: UNLEARNACC = 0
2: MAXEP = 100
3: EP = 0
4: $\bar{\mathcal{D}}_r \leftarrow$ Generate adversarial examples with Algorithm 1
5: $\bar{\Omega} \leftarrow$ Measure weight importance with Algorithm 2
6: $\tilde{\boldsymbol{\theta}} \leftarrow \boldsymbol{\theta}$ // Set the pretrained model for unlearning
7: **while** UNLEARNACC $\neq 100$ **do**
8:    Minimize Eqn (6) or (7)
9:    UNLEARNACC = GetAccuracy($\mathcal{D}_f, g_{\tilde{\boldsymbol{\theta}}}$)
10:   **if** EP > MAXEP **then**
11:      break
12:      EP += 1
13:   **end if**
14: **end while**
15: **return** $\tilde{\boldsymbol{\theta}}$

---

## B. Details for RAWP

Algorithm 5 presents the pseudo code of RAWP. As an alternative to minimizing cross-entropy loss, we employed adversarial weight perturbation (AWP) [43] for the proposed

---

**Algorithm 5** The pseudo-code of overall unlearning process the case of using $\mathcal{L}_{\text{UL}}^{\text{MS}}$ with RAWP.

---
1: UNLEARNACC = 100
2: MAXEP = 100
3: EP = 0
4: $\tilde{\boldsymbol{\theta}} \leftarrow \boldsymbol{\theta}$ // Set the pretrained model for unlearning
5: $\hat{\boldsymbol{\theta}} \leftarrow \boldsymbol{\theta}$ // Set the pretrained model for AWP
6: **while** UNLEARNACC $\neq 0$ **do**
7:    $\hat{\boldsymbol{\theta}}_p \leftarrow \hat{\boldsymbol{\theta}}$ // Set a proxy model
8:    Minimize $\mathcal{L}_{\text{UL}}^{\text{MS}}$ for $\hat{\boldsymbol{\theta}}_p$
9:    Get gradients and update $\hat{\boldsymbol{\theta}}_p$
10:   Get difference $\hat{\boldsymbol{\theta}}_d$ of each layer $i$ following
       : $\hat{\boldsymbol{\theta}}_{d,i} = ||\tilde{\boldsymbol{\theta}}_i||_2 / (||\tilde{\boldsymbol{\theta}}_i - \hat{\boldsymbol{\theta}}_{p,i}||_2 + \epsilon) * (\tilde{\boldsymbol{\theta}}_i - \hat{\boldsymbol{\theta}}_{p,i})$
11:   Perturb $\tilde{\boldsymbol{\theta}}$ of each layer $i$ following
       : $\tilde{\boldsymbol{\theta}}_i = \tilde{\boldsymbol{\theta}}_i + \gamma\hat{\boldsymbol{\theta}}_{d,i}$
12:   UNLEARNACC = GetAccuracy($\mathcal{D}_f, g_{\tilde{\boldsymbol{\theta}}}$)
13:   **if** EP > MAXEP **then**
14:      break
15:      EP += 1
16:   **end if**
17: **end while**
18: **return** $\hat{\boldsymbol{\theta}}$

---

unlearning. The main difference between AWP and cross-entropy loss is that AWP does not directly minimize loss through gradient descent, but instead introduces a worst-case weight perturbation with a relative perturbation size constrained by $\gamma$. In the manuscript, we reported experimental results using the default value of $\gamma = 0.01$.

Table 6 reports additional experiments results for various $\gamma$. We confirm that, for all $\gamma$ cases, RAWP achieves degraded performance especially for large $k$ sizes.

Table 6: Experimental results for various $\gamma$ of RAWP.

|   |   | $k = 16$ | $k = 64$ | $k = 128$ | $k = 256$ |
|---|---|---|---|---|---|
| RAWP ($\gamma = 0.001$) | $D_f$ | 0.0 | 0.0 | 0.0 | 0.0 |
|   | $D_r$ | 88.31 | 65.44 | 42.10 | 14.81 |
|   | $D_{test}$ | 81.66 | 60.90 | 39.16 | 14.34 |
| RAWP ($\gamma = 0.005$) | $D_f$ | 0.0 | 0.0 | 0.0 | 2.89 |
|   | $D_r$ | 81.70 | 55.09 | 28.60 | 13.45 |
|   | $D_{test}$ | 75.80 | 51.79 | 27.21 | 13.11 |
| RAWP ($\gamma = 0.01$) | $D_f$ | 0.0 | 0.0 | 4.06 | 6.17 |
|   | $D_r$ | 67.15 | 37.33 | 13.28 | 9.12 |
|   | $D_{test}$ | 63.10 | 35.47 | 12.91 | 9.12 |
| RAWP ($\gamma = 0.1$) | $D_f$ | 0.9 | 6.25 | 8.43 | 9.61 |
|   | $D_r$ | 10.97 | 9.16 | 10.00 | 10.00 |
|   | $D_{test}$ | 11.13 | 9.09 | 10.00 | 10.00 |

## C. Additional results on various models

**Results on various models.** Figure 5 shows unlearning results on CIFAR-100, but with different model architectures. We find that our methods effectively preserve knowledge outside the forgetting data, resulting in up to 40% boost in accu-
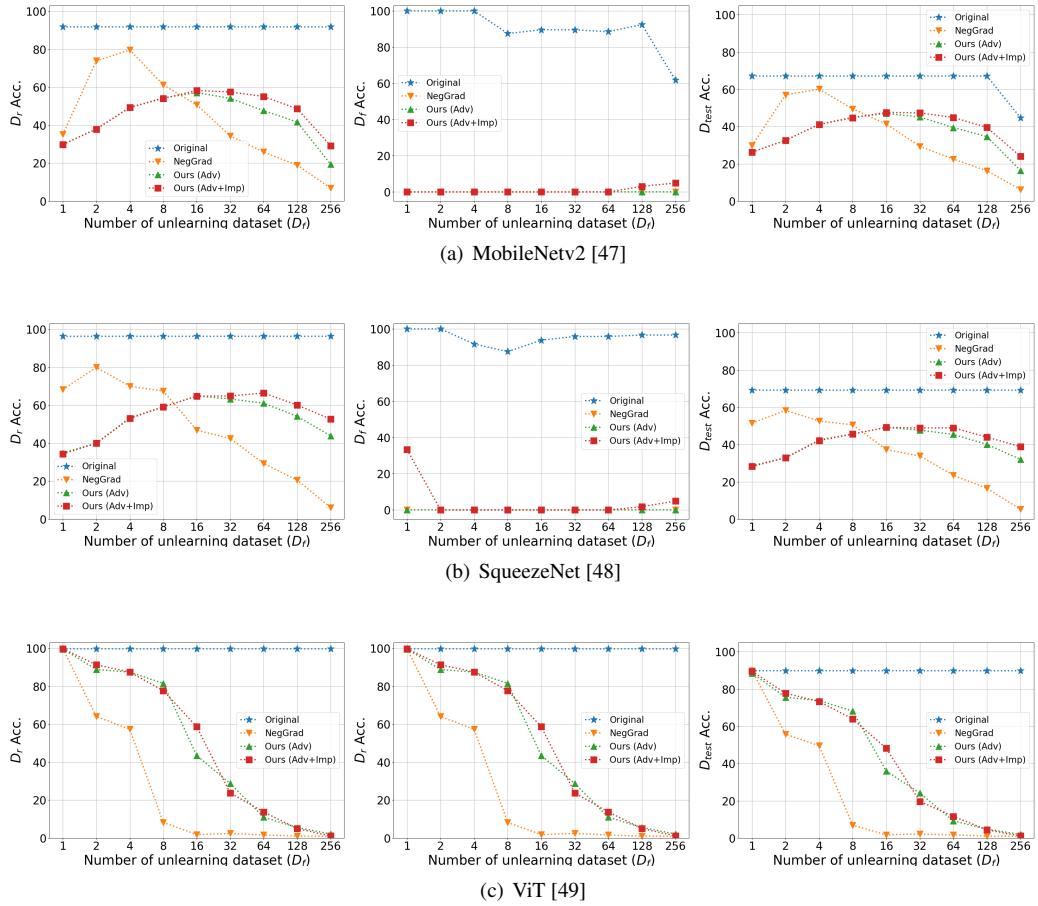
(a) MobileNetv2 [47]

(b) SqueezeNet [48]

(c) ViT [49]

Figure 5: Experimental results before and after unlearning varying $k$ instances from various models on CIFAR-100.

racy. NegGrad again outperforms our methods when $k = 4$, but soon breaks down when unlearning more instances. Interestingly, SqueezeNet and MobileNetv2 suffer from larger forgetting in $\mathcal{D}_r$ and $\mathcal{D}_{test}$ than ResNet-50, possibly due to the width being narrower as previously investigated by [51]. ViT also suffers from large forgetting, an observation consistent with previous work which showed that ViT suffers more catastrophic forgetting compared to other CNN-based methods in continual learning due to Transformer architectures requiring large amounts of data. We also evaluate the results of unlearning on ImageNet-1K with varying $k$ in Figure 6. Our proposed methods prevent forgetting knowledge about the rest data $\mathcal{D}_r$ better than NegGrad in all cases where k is greater than 8. At the same time, the methods effectively delete information about $\mathcal{D}_f$.

## References

[1] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.

[2] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15262–15271, 2021.

[3] Pierre Lison, Ildikó Pilán, David Sánchez, Montserrat Batet, and Lilja Øvrelid. Anonymisation models for text data: State of the art, challenges and future directions. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4188–4203, 2021.

[4] Rich Zemel, Yu Wu, Kevin Swersky, Toni Pitassi, and Cynthia Dwork. Learning fair representations. In *International conference on machine learning*, pages 325–333. PMLR, 2013.

[5] Preethi Lahoti, Alex Beutel, Jilin Chen, Kang Lee, Flavien Prost, Nithum Thain, Xuezhi Wang, and Ed Chi. Fairness without demographics through adversarially reweighted learning. *Advances in neural information processing systems*, 33:728–740, 2020.

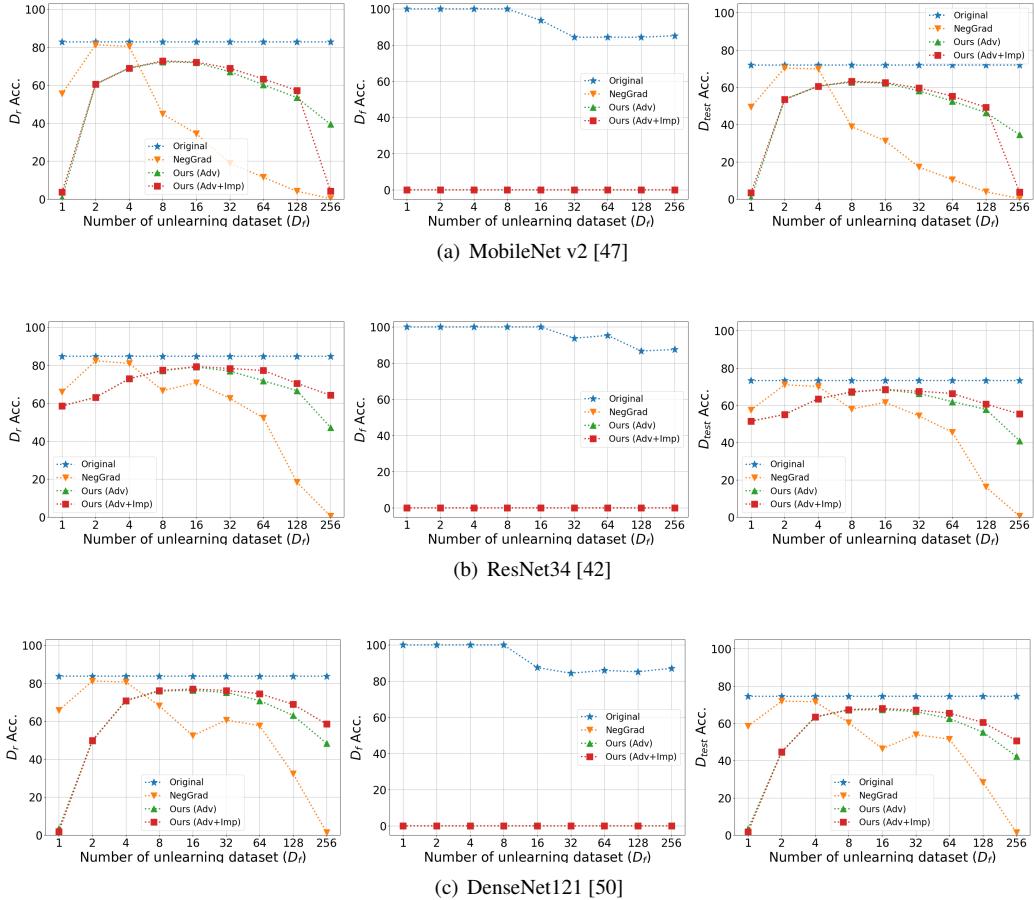[6] Melissa Heikkilä. What does gpt-3 "know" about me?, Aug 2022.

(a) MobileNet v2 [47]



(b) ResNet34 [42]



(c) DenseNet121 [50]

Figure 6: Experimental results before and after unlearning varying $k$ instances from various models on ImageNet-1K.

[7] Jeffrey Rosen. The right to be forgotten. *Stan. L. Rev. Online*, 64:88, 2011.

[8] Eduard Fosch Villaronga, Peter Kieseberg, and Tiffany Li. Humans forget, machines remember: Artificial intelligence and the right to be forgotten. *Computer Law & Security Review*, 34(2):304–313, 2018.

[9] Ananth Mahadevan and Michael Mathioudakis. Certifiable machine unlearning for linear models. *arXiv preprint arXiv:2106.15093*, 2021.

[10] 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.

[11] Jonathan Brophy and Daniel Lowd. Machine unlearning for random forests. In *International Conference on Machine Learning*, pages 1092–1104. PMLR, 2021.

[12] Ayush K Tarun, Vikram S Chundawat, Murari Mandal, and Mohan Kankanhalli. Fast yet effective machine unlearning. *arXiv preprint arXiv:2111.08947*, 2021.

[13] Vikram S Chundawat, Ayush K Tarun, Murari Mandal, and Mohan Kankanhalli. Zero-shot machine unlearning. *arXiv preprint arXiv:2201.05629*, 2022.

[14] Jingwen Ye, Yifang Fu, Jie Song, Xingyi Yang, Songhua Liu, Xin Jin, Mingli Song, and Xinchao Wang. Learning with recoverable forgetting. *arXiv preprint arXiv:2207.08224*, 2022.

[15] Youngsik Yoon, Jinhwan Nam, Hyojeong Yun, Dongwoo Kim, and Jungseul Ok. Few-shot unlearning by model inversion. *arXiv preprint arXiv:2205.15567*, 2022.

[16] 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*, pages 9304–9312, 2020.

[17] Junyaup Kim and Simon S Woo. Efficient two-stage model retraining for machine unlearning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4361–4369, 2022.

[18] Ronak Mehta, Sourav Pal, Vikas Singh, and Sathya N Ravi. Deep unlearning via randomized conditionally independent hessians. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10422–10431, 2022.

[19] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A survey on bias and fairness in machine learning. *ACM Computing Surveys (CSUR)*, 54(6):1–35, 2021.

[20] Laura Graves, Vineel Nagisetty, and Vijay Ganesh. Amnesiac machine learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 11516–11524, 2021.

[21] Sayna Ebrahimi, Franziska Meier, Roberto Calandra, Trevor Darrell, and Marcus Rohrbach. Adversarial continual learning. In *European Conference on Computer Vision*, pages 386–402. Springer, 2020.

[22] Rahaf Aljundi, Francesca Babiloni, Mohamed Elhoseiny, Marcus Rohrbach, and Tinne Tuytelaars. Memory aware synapses: Learning what (not) to forget. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 139–154, 2018.

[23] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.

[24] 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*, pages 248–255. Ieee, 2009.

[25] Yinzhi Cao and Junfeng Yang. Towards making systems forget with machine unlearning. In *2015 IEEE Symposium on Security and Privacy*, pages 463–480. IEEE, 2015.

[26] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.

[27] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.

[28] Alexey Kurakin, Ian Goodfellow, Samy Bengio, et al. Adversarial examples in the physical world, 2016.

[29] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.

[30] Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *2017 ieee symposium on security and privacy (sp)*, pages 39–57. IEEE, 2017.

[31] Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry. Adversarial examples are not bugs, they are features. *Advances in neural information processing systems*, 32, 2019.

[32] Pavlo Molchanov, Arun Mallya, Stephen Tyree, Iuri Frosio, and Jan Kautz. Importance estimation for neural network pruning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11264–11272, 2019.

[33] Zhuang Liu, Jianguo Li, Zhiqiang Shen, Gao Huang, Shoumeng Yan, and Changshui Zhang. Learning efficient convolutional networks through network slimming. In *Proceedings of the IEEE international conference on computer vision*, pages 2736–2744, 2017.

[34] Wei Wen, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. Learning structured sparsity in deep neural networks. *Advances in neural information processing systems*, 29, 2016.

[35] Jose M Alvarez and Mathieu Salzmann. Learning the number of neurons in deep networks. *Advances in neural information processing systems*, 29, 2016.

[36] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets. *arXiv preprint arXiv:1608.08710*, 2016.

[37] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.

[38] Arslan Chaudhry, Puneet K Dokania, Thalaiyasingam Ajanthan, and Philip HS Torr. Riemannian walk for incremental learning: Understanding forgetting and intransigence. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 532–547, 2018.

[39] Sangwon Jung, Hongjoon Ahn, Sungmin Cha, and Taesup Moon. Continual learning with node-importance based adaptive group sparse regularization. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 33, pages 3647–3658. Curran Associates, Inc., 2020.

[40] Rahaf Aljundi, Marcus Rohrbach, and Tinne Tuytelaars. Selfless sequential learning. In *International Conference on Learning Representations (ICLR)*, 2019.

[41] Sungmin Cha, Hsiang Hsu, Taebaek Hwang, Flavio Calmon, and Taesup Moon. Cpr: Classifier-projection regularization for continual learning. In *International Conference on Learning Representations*.

[42] 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*, pages 770–778, 2016.

[43] Dongxian Wu, Shu-Tao Xia, and Yisen Wang. Adversarial weight perturbation helps robust generalization. *Advances in Neural Information Processing Systems*, 33:2958–2969, 2020.

[44] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008.

[45] Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural network representations revisited. In *International Conference on Machine Learning*, pages 3519–3529. PMLR, 2019.

[46] Mariya Toneva, Alessandro Sordoni, Remi Tachet des Combes, Adam Trischler, Yoshua Bengio, and Geoffrey J Gordon. An empirical study of example forgetting during deep neural network learning. *arXiv preprint arXiv:1812.05159*, 2018.

[47] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4510–4520, 2018.

[48] Forrest N Iandola, Song Han, Matthew W Moskewicz, Khalid Ashraf, William J Dally, and Kurt Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and¡ 0.5 mb model size. *arXiv preprint arXiv:1602.07360*, 2016.

[49] 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.

[50] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017.

[51] Seyed Iman Mirzadeh, Arslan Chaudhry, Dong Yin, Timothy Nguyen, Razvan Pascanu, Dilan Gorur, and Mehrdad Farajtabar. Architecture matters in continual learning. *arXiv preprint arXiv:2202.00275*, 2022.