

# Layer Attack Unlearning: Fast and Accurate Machine Unlearning via Layer Level Attack and Knowledge Distillation

Hyunjune Kim<sup>1</sup>, Sangyong Lee<sup>2</sup>, Simon S. Woo<sup>1 2 3</sup>

<sup>1</sup> Department of Applied Data Science, Sungkyunkwan University, South Korea

<sup>2</sup> Department of Artificial Intelligence, Sungkyunkwan University, South Korea

<sup>3</sup> Department of Computer Science and Engineering, Sungkyunkwan University, South Korea  
hyunjune.kim@g.skku.edu, sang8961@g.skku.edu, swoo@g.skku.edu

## Abstract

Recently, serious concerns have been raised about the privacy issues related to training datasets in machine learning algorithms when including personal data. Various regulations in different countries, including the GDPR, grant individuals to have personal data erased, known as ‘the right to be forgotten’ or ‘the right to erasure’. However, there has been less research on effectively and practically deleting the requested personal data from the training set while not jeopardizing the overall machine learning performance. In this work, we propose a fast and novel machine unlearning paradigm at the layer level called layer attack unlearning, which is highly accurate and fast compared to existing machine unlearning algorithms. We introduce the Partial-PGD algorithm to locate the samples to forget efficiently. In addition, we only use the last layer of the model inspired by the Forward-Forward algorithm for unlearning process. Lastly, we use Knowledge Distillation (KD) to reliably learn the decision boundaries from the teacher using soft label information to improve accuracy performance. We conducted extensive experiments with SOTA machine unlearning models and demonstrated the effectiveness of our approach for accuracy and end-to-end unlearning performance.

## 1 Introduction

Deep neural networks (DNNs) have achieved significant progress and dramatic performance gains in challenging machine learning tasks in recent years. Among others, large amounts of available training datasets have been the foundation for enabling the revolution of large-scale models. However, recently, due to the privacy concerns raised by ChatGPT (Bourtoule et al. 2021; Burgess 2023), the training dataset would contain personal information or possible information that can leak one’s privacy. For example, many vision-based applications would involve training one’s face images, which are personally identifiable information (PII). Several nations have implemented some types of regulations, such as the General Data Protection Regulation (GDPR) (Mantelero 2013) and the EU/US Copyright Law (Kaye 2023; Kublik 2023), in order to address the potential misuse of personal information and further grant individuals the right to have personal data erased, known as ‘the right to be forgotten’ or ‘the right to erasure.’ The goal

of such regulations is to provide data owners the right to request and erase the personal or copyrighted data they want if they have not agreed and consented in the first place.

Therefore, companies using personal data should delete the requested data from the training set. One potential approach for corporations to mitigate the aforementioned concerns involves the exclusion of the required dataset from the training dataset, followed by a complete retraining process from scratch. Nevertheless, as models like ChatGPT get bigger and datasets grow, retraining them from scratch requires excessive computational resources and time.

*Machine unlearning* has emerged to tackle this challenge, allowing ML models to discard specific data selectively. (Bourtoule et al. 2021) Machine unlearning can be divided into two primary strategies: *instance-wise* and *class-wise* unlearning. The former involves forgetting knowledge related to specific instances from ML models, while the latter, which we focus on, completely removes particular classes from ML models. For example, face recognition and social media classification systems may need to erase data related to specific religion, nationality, age, disease, gender, etc., for security and privacy reasons. A few approaches (Chen et al. 2023; Cha et al. 2023) have explored the adversarial attacks for unlearning by harnessing the forgetting data’s noise to navigate the adjacent latent space. However, they used the original PGD (Madry et al. 2017) for unlearning, which can be slow.

In this work, we propose **Layer Attack Unlearning**, a fast and novel machine unlearning algorithm to tackle the class-wise unlearning problem. Our approach first introduces **Partial-PGD**, which is a new adversarial attack generation strategy to efficiently search the close vicinity of the data points to delete (See Fig. 1). Our proposed Partial-PGD is designed only to attack fully connected (classification) layer for probing the neighboring latent space to shift the forgetting data. Surprisingly, we do not utilize any feature layer information while achieving efficiency and accuracy. As shown in Fig. 1, Partial-PGD is much more efficient than the original PGD, as it can create adversarial examples only via the classification layer.

In particular, Hinton (2022)’s Forward-Forward (FF) algorithm has inspired us, and we provide the foundation of the concept of layer-based attack for machine unlearning based on FF. According to Hinton (2022), each layer un-

dergoes individualized training in the Forward-Forward algorithm to achieve its specific objectives. Similarly, in line with FF research, we aim to accomplish machine unlearning objectives at layers with characteristics directly relevant to data and features we want to forget. Hence, we focus on performing **machine unlearning at the layer level** rather than using the entire model. Our layer-wise unlearning approach clearly avoids unnecessary loss calculations during the unlearning process. Furthermore, updating only the layers’ weights related to forgetting data will ensure a reduction in computational costs.

Finally, we employ Knowledge Distillation (KD) (Hinton, Vinyals, and Dean 2015) to modify the decision boundary for the forgetting data and preserve the decision boundary for the retain data. The primary objective in unlearning is to utilize hard labels and acquire soft label information from the teacher model for unlearning tasks to maintain performance. We show that it achieves a stable placement of forgetting data in the space subjected to carefully created adversarial examples. We incorporated KD into our final loss function to improve performance.

Our main contributions are summarized as follows:

- We introduce Layer Attack Unlearning (LAU) algorithm, which is a novel and fast unlearning method by proposing Partial-PGD and performing unlearning at the layer level.
- In addition, we propose KD method to further improve the overall accuracy and data erasure performance by effectively distilling the decision boundary knowledge from the teacher model for unlearning task.
- Our extensive experimental results with seven baselines with four different backbones, including ViT over three other datasets, show that our approach outperforms previous SOTA methods in accuracy and time performance while completely forgetting the requested class.

## 2 Related Work

There are two main approaches to the current machine unlearning problem in DNNs. The first involves considering unlearning during the learning process, while the second focuses on fine-tuning. This paper will refer to the approach that considers the learning process as “data-driven” and the approach that involves fine-tuning as “model-agnostic.”

### 2.1 Data-Driven Unlearning Methods

A “data-driven” approach utilizes data-centric strategies such as partitioning and augmentation (Nguyen et al. 2022) to address unlearning. SISA (Bourtoule et al. 2021) and Selective Forgetting (Shibata et al. 2021) are two representative data-driven unlearning methods. In SISA, data is divided into shard units, sequentially trained in slices, and multiple model checkpoints are created. Once an unlearning query is requested, it reverts the query to the checkpoint before learning and retrains this reverted query with the ensemble technique. However, it is challenging to calculate the probability of encountering unlearning queries on data points.

On the other hand, Selective Forgetting (Shibata et al. 2021) involves lifelong learning to perform unlearning. A “mnemonic code” signal is embedded in the data during

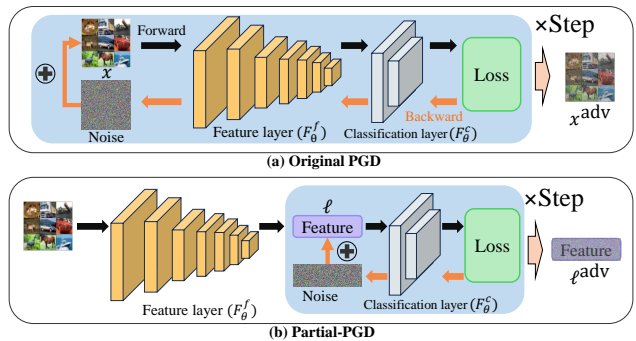


Figure 1: Illustration of the original PGD vs. Partial-PGD. While the original PGD involves backpropagation to compute  $x^{\text{adv}}$  with respect to input  $x$  for all the layers, (b) Partial-PGD computes  $\ell^{\text{adv}}$  in  $\mathcal{F}_\theta^c$  after passing  $x$  through  $\mathcal{F}_\theta^f$  to calculate  $\ell$ . Step in both (a) and (b) indicates the iteration.

training. During the unlearning process, the mnemonic code information is selectively incorporated into the loss function to remove forgetting data. This strategy requires storing mnemonic codes for all data points, considering unlearning queries before building the original model. This could be more practical in a real-world scenario.

### 2.2 Model-Agnostic Unlearning Methods

A “model-agnostic” approach is a methodology for handling the unlearning process by adjusting the model’s learning parameters to achieve data unlearning (Nguyen et al. 2022). Such approaches include various methods such as Summation form (Cao and Yang 2015), Negative Gradient (Golatkar, Achille, and Soatto 2020), Fisher Forgetting (Golatkar, Achille, and Soatto 2020), Boundary unlearning (Chen et al. 2023), Instance-wise Unlearning (Cha et al. 2023), etc. Some methods utilize adversarial attacks to the original model to avoid naively excluding and deleting forgetting data. Among the mentioned algorithms, approaches like ours include Boundary unlearning and Instance-wise Unlearning. These two algorithms perform unlearning by utilizing adversarial attacks to transition forgettable data to nearby spaces. However, a significant difference between our approach and these methods lies in the target of the attack. Our approach directs the unlearning process towards layers with specific classification objectives instead of using entire layers. Furthermore, we aim to introduce effective ways of utilizing PGD in unlearning.

## 3 Our Approach

The main objective of our approach is to accurately and efficiently perform *class-wise* unlearning, which is to completely remove specific classes from the classification model. In this section, we describe our Partial-PGD, KD architecture, and our connection to the FF algorithm.

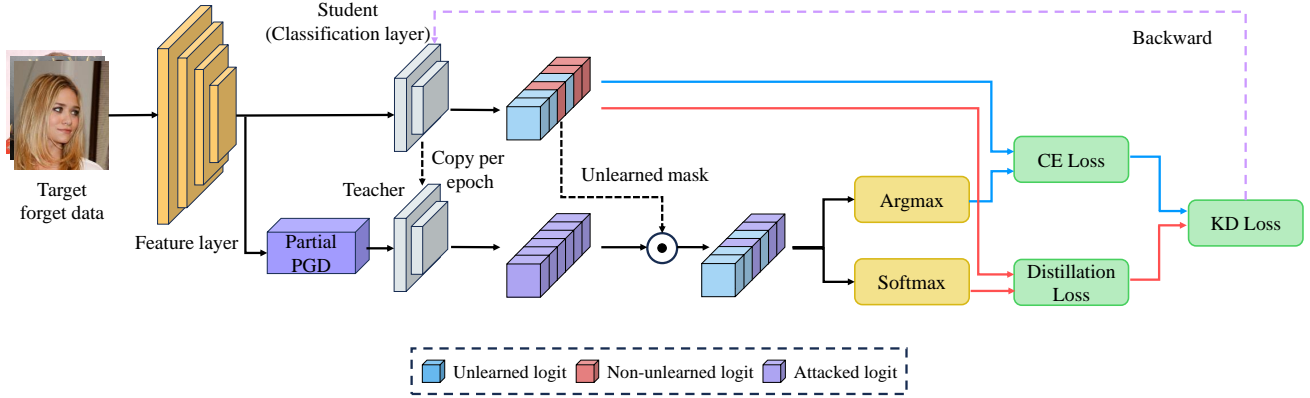


Figure 2: The overall procedure of our approach. Our method involves the unlearning task on the classification layer instead of the entire model, where each classification layer represents the student and the teacher model. For the unlearning task, we perform Knowledge Distillation by combining the teacher logit and student logit via the unlearned mask. The teacher logit is derived from the adversarial examples obtained after applying Partial-PGD.

### 3.1 Preliminaries and Notations

First, we formulate a machine unlearning problem as follows: We define a training dataset  $D_{\text{train}} = \{x^i, y^i\}_{i=1}^N$ , consisting of inputs  $x^i \in X$  and their corresponding class labels  $y^i \in Y$ . The forgetting dataset  $D_f$  is a subset of  $D_{\text{train}}$  that we intend to forget from the pre-trained model. Conversely, the retain dataset  $D_r = D_{\text{train}} \setminus D_f$  is the dataset we want to preserve the overall performance.

Next, we define the original model  $\mathcal{M}_\theta : \mathbb{R}^n \rightarrow \mathbb{R}^n$ , which comprises a set of feature layers denoted by  $\mathcal{F}_\theta^f : \mathbb{R}^n \rightarrow \mathbb{R}^n$  and a fully connected layer denoted as  $\mathcal{F}_\theta^c : \mathbb{R}^n \rightarrow \mathbb{R}^n$ , where  $\theta$  represents the optimal parameters for the model trained on  $D_{\text{train}}$ . The following provides a compositional representation of the model  $\mathcal{M}_\theta$  as  $\mathcal{F}_\theta^c \circ \mathcal{F}_\theta^f$ . Also, we denote  $x^{\text{adv}}$  to represent the adversarial examples (Goodfellow, Shlens, and Szegedy 2014) for the input data  $x$ . In particular, we define  $\ell^{\text{adv}}$  as the adversarial example from Partial-PGD, generated from the intermediate latent feature  $\ell$  obtained from the outputs of  $\mathcal{F}_\theta^f$ , as shown in Fig. 1.

### 3.2 Partial-Projected Gradient Descent (PGD)

The main reason for employing adversarial examples is to search and identify neighboring candidate spaces more effectively that will assign the forgetting data samples. Assigning forgetting classes to random or irrelevant classes can dramatically reduce downstream task performance.

Therefore, carefully exploring the neighboring space allows us not only to forget  $D_f$  but also to preserve the decision boundary of other classes. Hence, adversarial attacks (Madry et al. 2017; Chen et al. 2023) can be explored below:

$$x^{t+1} = \Pi(x^t + (\epsilon \cdot \text{sign}(\nabla_x \mathcal{L}(x, y, \theta))), \quad (1)$$

where the parameter  $\theta$  represents the weights of the target

model under attack, and generated noise for crafting adversarial examples is produced by computing the gradient  $\nabla_x \mathcal{L}$  of the loss function  $\mathcal{L}$  with respect to the input  $x$ . This noise is added to  $x^t$  and then projected using the projection method  $\Pi$  to calculate  $x^{t+1}$ , which is repeated  $t$  times. Once  $x^{t+1}$  represents  $x^{\text{adv}}$ , it is an adversarial example.

However, we clarify the purpose of adversarial examples used in our work, which differs from prior approaches. The original PGD approach may generate excessive noise and slow the unlearning process considerably. Therefore, there is no need to calculate gradients throughout the entire model to create adversarial examples.

Hence, our proposed Partial-PGD utilizes  $\mathcal{F}_\theta^c$  to generate adversarial examples for the unlearning process, as shown in Fig. 1. This technique effectively identifies the neighboring space to allocate  $D_f$ , the forgetting data, similar to conventional PGD. However, it significantly reduces unlearning time by omitting feature layer information, as depicted in Fig. 1. We define our Partial-PGD as follows:

$$\ell^{t+1} = \Pi(\ell^t + (\epsilon \cdot \text{sign}(\nabla_\ell \mathcal{L}(\ell, y, \theta))), \quad (2)$$

where Partial-PGD applies an adversarial attack to the intermediate latent  $\ell$  obtained from  $\mathcal{F}_\theta^f$ , where  $\ell$  undergoes gradient computation based solely on passing through  $\mathcal{F}_\theta^c$ . Then, the result is mapped to the nearby space of a different label of  $\ell$  and becomes  $\ell^{\text{adv}}$ , which we use for unlearning  $D_f$  as knowledge to be forgotten.

### 3.3 Layer Unlearning

While other approaches use entire layers for unlearning, we focus on unlearning only the relevant layers. Inspired by the FF technique, we focus on the classification layer  $\mathcal{F}_\theta^c$  to forget specific classes in the model for class-wise unlearning. Therefore, our layer unlearning focuses on only modifying the parameters of  $\mathcal{F}_\theta^c$  tied to classification instead of the entire layers and model  $\mathcal{M}_\theta$  to forget  $D_f$  effectively.

We define the following equation to describe our unlearning process, where we focus on the  $\mathcal{F}_\theta^c$  during the unlearning process to remove  $D_f$  from the model:

$$\mathcal{M}_{\theta^*} = \mathcal{F}_{\theta^*}^c \circ \mathcal{F}_\theta^f, \quad (3)$$

where  $\theta^*$  is the ideal parameters after forgetting  $D_f$ .

We show that layer unlearning accelerates the unlearning process by selectively updating relevant layer weights and optimizing efficiency. Interestingly, it outperforms models with whole layers in accuracy.

### 3.4 End-to-End Unlearning Process

We describe our end-to-end unlearning process, where we apply the KD to improve the overall performance further. As illustrated in Fig. 2, the classification layer  $\mathcal{F}_\theta^c$  serves as our *student* model  $\mathcal{S}_\theta$ . Additionally, at the beginning of each epoch, we duplicate the  $\mathcal{S}_\theta$  as our *teacher*  $\mathcal{T}_\theta$ . The model uses forgetting data  $D_f$  as input to create an intermediate latent feature  $\ell_f$  through the feature layer  $\mathcal{F}_\theta^f$ . Then,  $\ell_f$  becomes an adversarial example  $\ell_f^{\text{adv}}$  after applying a Partial-PGD on the  $\mathcal{T}_\theta$ .

Next,  $\ell_f$  and  $\ell_f^{\text{adv}}$  are passed through  $\mathcal{S}_\theta$  and  $\mathcal{T}_\theta$ , respectively, becoming logits for each student and teacher, as shown in Fig. 2. Then, the logit obtained from  $\mathcal{S}_\theta$  is compared with the ground truth  $y_f$ . If a discrepancy is observed, it is considered unlearned. Then, the unlearned logit replaces the adversarial logit from  $\mathcal{T}_\theta$ . This student’s logit is used to compute the cross-entropy loss as follows:

$$\mathcal{L}_{CE} = \begin{cases} CE(\mathcal{S}_\theta(\ell_f), y_f^{\text{adv}}) & \text{if } y_{\mathcal{S}_\theta} = y_f \\ CE(\mathcal{S}_\theta(\ell_f), y_{\mathcal{S}_\theta}) & \text{otherwise,} \end{cases} \quad (4)$$

where  $y_{\mathcal{S}_\theta}$  represents the predicted label from  $\mathcal{S}_\theta(\ell_f)$ , and CE is the cross-entropy function. This loss leaves the unlearned data in a state, where it makes wrong (unlearned) predictions. If not, it is trained to be a predicted label  $y_f^{\text{adv}}$  of adversarial logit, leading to its unlearning process. Next, let  $Z$  be the double Softmax representation, which is defined as:

$$Z = \begin{cases} \sigma(\mathcal{T}_\theta(\ell_f^{\text{adv}})) & \text{if } y_{\mathcal{S}_\theta} = y_f \\ \sigma(\mathcal{S}_\theta(\ell_f)) & \text{otherwise,} \end{cases} \quad (5)$$

where  $\sigma$  represents Softmax function. In Eq. 5, we performed double Softmax to distill knowledge by adjusting the probability distribution of the output from  $\mathcal{T}_\theta$ . This approach is intended to convey soft label information to  $\mathcal{S}_\theta$ . Exclusively unlearning  $\mathcal{F}_\theta^c$  maintains the decision boundaries of retain data, and slightly improves the overall accuracy. But, layer unlearning without double Softmax showed variable accuracy, as shown in the Fashion-MNIST dataset (Xiao, Rasul, and Vollgraf 2017). We show this effect in Section 4.3. Next, we define our distillation loss as follows:

$$\mathcal{L}_{DI} = \text{KL} \left( \sigma \left( \frac{\mathcal{S}_\theta(\ell_f)}{T} \right), \sigma \left( \frac{Z}{T} \right) \right), \quad (6)$$

where knowledge is distilled from  $Z$  of  $\mathcal{T}_\theta$  and KL is the KL divergence. During distillation, the computation of loss

### Algorithm 1: End-to-End Unlearning Process

---

**Input:**  $\mathcal{F}_\theta^f, \mathcal{F}_\theta^c, D_f$   
**Parameter:** Learning rate  $\eta$ , Hyper-parameters  $\alpha$ , Temperature  $T$ , Number of Epochs  $E$   
**Output:**  $\mathcal{M}_{\theta^*}$

- 1:  $\mathcal{S}_\theta \leftarrow \mathcal{F}_\theta^c$
- 2:  $\theta^* \leftarrow \theta$
- 3: **for**  $e$  in range  $E$  **do**
- 4:    $\mathcal{T}_{\theta^*} \leftarrow \mathcal{S}_{\theta^*}$
- 5:    $\mathcal{L} \leftarrow (1 - \alpha) \cdot \mathcal{L}_{CE} + \alpha \cdot T^2 \cdot \mathcal{L}_{DI}$
- 6:    $\theta^* \leftarrow \theta^* - \eta \cdot \mathcal{L}$
- 7:   **if**  $\mathcal{F}_{\theta^*}^c \circ \mathcal{F}_\theta^f(X_f) \neq Y_f$  **then**
- 8:     **break**
- 9:   **end if**
- 10: **end for**
- 11:  $\mathcal{M}_{\theta^*} \leftarrow \mathcal{F}_{\theta^*}^c \circ \mathcal{F}_\theta^f$
- 12: **return**  $\mathcal{M}_{\theta^*}$

---

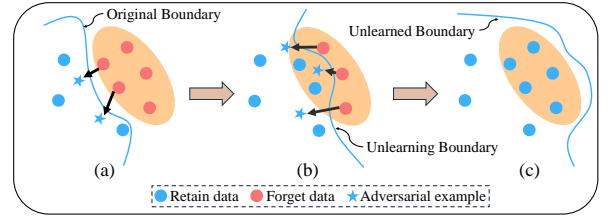


Figure 3: Boundary evolution in the unlearning process. As shown in (a), the original model receives the initial knowledge about the boundary. As the epoch progresses, the boundary information updates as depicted in (b) and (c) from the distilled knowledge.

$\mathcal{L}_{DI}$  between the outputs of  $\mathcal{S}_\theta$  and  $\mathcal{T}_\theta$  focuses on creating a similar boundary to the teacher model, ensuring performance while removing information of  $D_f$ . The temperature  $T$  is a hyper-parameter. Generally, increasing  $T$  will generate smoother soft labels that assists  $\mathcal{S}_\theta$  in mimicking  $\mathcal{T}_\theta$ . The effects of changes in  $T$  are described in Suppl. Mat.

Using  $\mathcal{L}_{CE}$  and  $\mathcal{L}_{DI}$ , our final loss function is constructed as follows:

$$\mathcal{L} = (1 - \alpha) \cdot \mathcal{L}_{CE} + \alpha \cdot T^2 \cdot \mathcal{L}_{DI}, \quad (7)$$

where the value of  $\alpha$  represents the weight assigned to the loss between  $\mathcal{L}_{CE}$  and  $\mathcal{L}_{DI}$ . As a hyper-parameter,  $\alpha$  ranges from 0 to 1. Assigning additional weight to  $\mathcal{L}_{CE}$  may boost unlearning time but decrease performance. Conversely, if we provide more weight to  $\mathcal{L}_{DI}$ , the unlearning speed may slow down but can increase accuracy. We conducted the ablation study for  $\alpha$  values to capture the trade-off. The effects of changes in the exponent of  $T^2$  are described in Suppl. Mat.

In addition, we provide the end-to-end unlearning process in Alg. 1. We distill knowledge from  $\mathcal{T}_\theta$ , while gradually reducing boundaries. Algorithm 1 finishes either when all epochs are completed or when  $D_f$  becomes unlearned within a batch during an epoch. Finally, we obtain our unlearning model  $\mathcal{M}_{\theta^*}$  by combining  $\mathcal{F}_\theta^f$  with the classification layer,  $\mathcal{F}_{\theta^*}^c$ , as shown in Eq. 3.

**Summary.** In Fig. 3, we pictorially describe our end-to-end unlearning process by displaying the boundary change for the retain and forgetting data.

## 4 Experimental Results

We experiment and evaluate popular unlearning benchmarks used in other unlearning research (Golatkar, Achille, and Soatto 2020; Chen et al. 2023; Cha et al. 2023) on image classification tasks.

**Datasets and Models.** We conducted experiments on CIFAR-10 (Krizhevsky, Hinton et al. 2009), Fashion-MNIST (Xiao, Rasul, and Vollgraf 2017), and VGGFace2 (Cao et al. 2018) datasets. For the VGGFace2 dataset, we randomly select ten individuals from a training dataset containing over 600 images, ensuring a balanced gender distribution. Furthermore, we perform training from scratch for three different architectures: VGG16 (Simonyan and Zisserman 2014), ResNets (He et al. 2016), and ViT (Dosovitskiy et al. 2020).

**Baseline Approaches.** The subsequent unlearning baseline methods are used:

- 1) **Original:** We train the model on the  $D_{\text{train}}$  dataset before undergoing the unlearning process.
- 2) **Retrain:** We train the model from scratch utilizing  $D_r$  as the retrained model, an optimal unlearning strategy.
- 3) **Negative Gradient (NG)** (Golatkar, Achille, and Soatto 2020): We fine-tune the **Original** with  $D_f$  by following the direction of gradient ascent.
- 4) **Fine-tune** (Golatkar, Achille, and Soatto 2020): We fine-tune the **Original** using  $D_r$  with a large learning rate.
- 5) **Random Label** (Golatkar, Achille, and Soatto 2020): We fine-tune the **Original** by assigning arbitrary labels randomly to  $D_f$ .
- 6) **Fisher Forgetting** (Golatkar, Achille, and Soatto 2020): The Fisher Forgetting model identifies influential parameters significantly affecting  $D_f$  and then introduces noise to neutralize their impact.
- 7) **Boundary Shrink** (Chen et al. 2023): We create adversarial examples from  $D_f$  and assign new adversarial labels to shrink towards different classes.
- 8) **IWU** (Cha et al. 2023): Generating adversarial instances for distinct labels via  $D_f$  and incorporating a regularization term. While initially designed for instance-wise learning, we adapt this method for class-wise unlearning problems.

**Implementation Details and Evaluation Metrics.** Our method and other baselines are implemented in Python 3.7 and use the PyTorch library (Paszke et al. 2019), employing a single NVIDIA GeForce RTX 3090 GPU. The initial model was trained using an LR scheduler and an SGD optimizer with specific settings (momentum: 0.9, weight decay:  $5 \times 10^{-4}$ , initial learning rate: 0.01). For the unlearning phase, we employ the SGD optimizer and conduct experiments with varying learning rates (ranging from 0.001 to 0.01), KD  $\alpha$  values (ranging from 0.3 to 0.7), KD temperature  $T$  value (fixed at 4), and Partial-PGD values (ranging from 0.4 to 1.0). As defined,  $D_f$  and  $D_r$  represent the forgetting and retain data, respectively. Additionally,  $D_{t_f}$  corresponds to the test forgetting data, and  $D_{t_r}$  represents the

test retain data. We assess the accuracy of all four different datasets.

### 4.1 Accuracy Performance

To achieve the best unlearning performance, it should completely forget information related to  $D_f$ . Therefore, guaranteeing accuracy on a par with those achieved by the **Retrain** for both  $D_f$  and  $D_r$  will be the best. Table 1 presents test results from different classification models, datasets, and metrics. The tested models include VGG16, ResNet18, ResNet50, and ViT. The datasets used for testing were CIFAR-10, Fashion-MNIST, and VGGFace2. In addition to the accuracy metric, we evaluate the performance using the unlearning score (US) as follows:

$$\text{US}(\text{acc}_r, \text{acc}_f) = \frac{\exp(\frac{\text{acc}_r}{100}) + \exp(1 - \frac{\text{acc}_f}{100}) - 2}{2 \cdot (\exp(1) - 1)}, \quad (8)$$

where  $\text{acc}_r$  and  $\text{acc}_f$  denote the accuracy of the retain and forgetting dataset, respectively. If the  $D_{t_r}$  approaches 100% and  $D_{t_f}$  approaches 0%, the US metric approaches 1, indicating a stable result on the unlearning process. We provide a more detailed explanation of why this metric is useful for unlearning in Suppl. Mat.

Finally, Table 1 presents the performance of each unlearning method for a specific single class in the aforementioned datasets. We measure the accuracy for datasets  $D_r$ ,  $D_f$ ,  $D_{t_r}$ , and  $D_{t_f}$ , along with the metric US. For the **NG**, the unstable variability in the loss of negative gradient contributes to less favorable overall performance results. **Fine-tune** shows strong performance in forgetting and retaining information. Nevertheless, this methodology requires utilizing the complete dataset  $D_r$  during training. Such extensive data is time-consuming, and we analyze and compare their worse time performance in Table 2. In the case of **Random Label**, except for VGGFace2’s ResNet18, most cases have poor accuracy and US. Due to the random nature of forgetting, converging towards arbitrary labels in the classification space is challenging, resulting in low performance.

**Fisher Forgetting** exhibits poor performance, with low accuracy and US on the overall test. Also, the Fisher matrix information required a significant amount of time. For **Boundary Shrink**, they also utilized adversarial attack examples, but they used the hard label information of the attack examples on  $D_f$ , which resulted in an unstable unlearning process. **IWU** approach involves utilizing adversarial attack examples while incorporating regularization to achieve a stable unlearning process. However, this gains an average US of 0.8587 in the overall test.

Finally, **Ours** completely removes the forgetting dataset (0% accuracy) on all the test cases and retains the highest unlearning performance. The accuracy for both  $D_f$  and  $D_{t_f}$  reaches 0, while the accuracy for  $D_r$  and  $D_{t_r}$  is comparable to or sometimes even higher than the **Retrain**. Also, ours demonstrates superior performance compared to almost all baseline models across various scenarios, with a high US average of 0.9443. Our approach that utilizes Partial-PGD and KD-based unlearning processes on layers with explicit objectives clearly achieves the best unlearning performance.



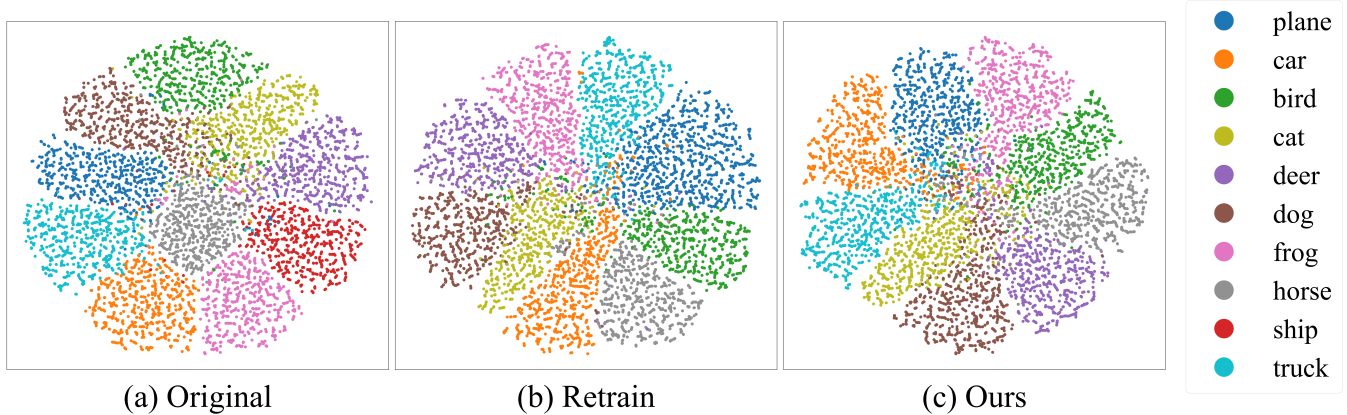


Figure 5: Visualization of decision boundary for the CIFAR-10 dataset in ResNet18, where each point represents a sample colored with the predicted classes. Red dots in (a) are the data to be removed, which are not showing in (b) and (c), indicating the successful unlearning. Similar plots for other models are provided in Suppl. Mat.

Table 4: Effect of Softmax vs. Double Softmax.

	w/o Double Softmax			w/ Double Softmax		
	$D_{tr}$	$D_{tf}$	Time (s)	$D_{tr}$	$D_{tf}$	Time (s)
VGG16	84.74	0	10.9	93.89	0	8.75
ResNet18	91.42	0.1	25.87	94.54	0	5.19
ResNet50	80.91	0	93.49	94.48	0	9.13
ViT	87.01	0	61.37	87.44	0	13.39

Table 5: The changes in time and accuracy performance with the reduction in  $D_f$  data on the CIFAR-10.

	Model	VGG16		ResNet18		ResNet50		ViT	
	Total Extra Data Used	2,500	500	2,500	500	2,500	500	2,500	500
Metrics	$D_{tr}$	92.42	92.38	93.51	93.38	93.63	93.37	81.14	81.6
	$D_{tf}$	0	0	0	0	0	0	0	0.1
	US	0.9422	0.9420	0.9503	0.9493	0.9512	0.9493	0.8640	0.8662
	Time	1.91	1.21	2.28	1.45	3.81	1.62	25.63	14.55

distillation loss. We have coined this method ‘‘Double Softmax’’, where Double Softmax enhances the robustness of our method across diverse datasets and models. And, Table 4 presents unlearning performance with and without double Softmax in our methods on the Fashion-MNIST dataset.

**Data Usage Ratio.** The class-specific  $D_f$  dataset for one class in CIFAR-10 contains 5,000 samples. As shown in Table 5, we reduced the dataset size to 50% (2,500) and 10% (500) for each model to perform the unlearning task. We measure the accuracy, US, and execution time of  $D_{tr}$ ,  $D_{tf}$ . In the following scenario, all models completed the unlearning for 2,500 samples, but ViT still had 0.1% retaining for 500 samples. The execution speed increases as the size of  $D_f$  decreases. Our experiment shows the potential for achieving superior unlearning performance by focusing on critical subsets of  $D_f$  rather than employing the complete dataset, saving time nearly seven times.

**Hyper-parameter  $\alpha$  in KD.** As shown in Fig. 4, we examine the accuracy variation of  $D_{tr}$  and  $D_{tf}$  with respect

to changes in the hyper-parameter  $\alpha$  in Eq. 6. As the  $\alpha$  approaches zero, it exclusively prioritizes the removal of  $D_f$  without taking into account any information from  $D_r$ . Consequently, the information about  $D_{tf}$  is completely removed, resulting in a decrease in the accuracy of  $D_{tr}$ . As  $\alpha$  approaches one, heavily relying on the teacher model for retaining information increases  $D_{tf}$  accuracy, indicating ineffective unlearning. Therefore, selecting the appropriate  $\alpha$  value can maximize unlearning performance. Consequently, we used  $\alpha$  ranging from 0.4 to 0.6 in this work. In more detail, the effects of changes in  $\alpha$  are described in Suppl. Mat.

#### 4.4 Visualization on Decision Boundary

Figure 5 presents the **Original**, **Retrain**, and **Ours** using t-SNE on the CIFAR-10 dataset. The red dots represent samples of ship images, indicated as  $D_f$ . As shown in Fig. 5(b),  $D_f$  was totally misclassified in the **Retrain**. On the other hand, Our unlearning method produces results correctly, as shown in Fig. 5(c), where the decision boundary of  $D_f$  has been successfully absorbed into the surrounding space.

## 5 Conclusion

In this paper, we introduce a novel and fast machine unlearning algorithm, layer attack unlearning, which is the new layer-based unlearning paradigm. Our work proposes Partial-PGD, layer unlearning method, and KD end-to-end framework to improve the overall accuracy performance while completely removing the forgetting dataset. Through extensive experimental evaluations, we demonstrated that modifying only specific layers’ learning objectives can lead to successful unlearning. Our approach effectively decreases both the number of parameters and their updates (computational cost), consequently reducing the overall time required for unlearning. We believe our layer attack unlearning paves a new way for future research in effectively addressing various unlearning challenges.

## 6 Acknowledgments

The authors would thank anonymous reviewers. Simon S. Woo is the corresponding author. This work was partly supported by Institute for Information & communication Technology Planning & evaluation (IITP) grants funded by the Korean government MSIT: (No. 2022-0-01199, Graduate School of Convergence Security at Sungkyunkwan University), (No. 2022-0-01045, Self-directed Multi-Modal Intelligence for solving unknown, open domain problems), (No. 2022-0-00688, AI Platform to Fully Adapt and Reflect Privacy-Policy Changes), (No. 2021-0-02068, Artificial Intelligence Innovation Hub), (No. 2019-0-00421, AI Graduate School Support Program at Sungkyunkwan University), and (No. RS-2023-00230337, Advanced and Proactive AI Platform Research and Development Against Malicious deepfakes).

## References

- Bourtole, L.; Chandrasekaran, V.; Choquette-Choo, C. A.; Jia, H.; Travers, A.; Zhang, B.; Lie, D.; and Papernot, N. 2021. Machine unlearning. In *2021 IEEE Symposium on Security and Privacy (SP)*, 141–159. IEEE.
- Burgess, M. 2023. ChatGPT Has a Big Privacy Problem. <https://www.wired.com/story/italy-ban-chatgpt-privacy-gdpr/>.
- Cao, Q.; Shen, L.; Xie, W.; Parkhi, O. M.; and Zisserman, A. 2018. Vggface2: A dataset for recognising faces across pose and age. In *2018 13th IEEE international conference on automatic face & gesture recognition (FG 2018)*, 67–74. IEEE.
- Cao, Y.; and Yang, J. 2015. Towards making systems forget with machine unlearning. In *2015 IEEE symposium on security and privacy*, 463–480. IEEE.
- Cha, S.; Cho, S.; Hwang, D.; Lee, H.; Moon, T.; and Lee, M. 2023. Learning to unlearn: Instance-wise unlearning for pre-trained classifiers. *arXiv preprint arXiv:2301.11578*.
- Chen, M.; Gao, W.; Liu, G.; Peng, K.; and Wang, C. 2023. Boundary Unlearning. *arXiv preprint arXiv:2303.11570*.
- Dosovitskiy, A.; Beyer, L.; Kolesnikov, A.; Weissenborn, D.; Zhai, X.; and Unterthiner, T. 2020. Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.
- Golatkar, A.; Achille, A.; and Soatto, S. 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.
- Goodfellow, I. J.; Shlens, J.; and Szegedy, C. 2014. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 770–778.
- Hinton, G. 2022. The forward-forward algorithm: Some preliminary investigations. *arXiv preprint arXiv:2212.13345*.
- Hinton, G.; Vinyals, O.; and Dean, J. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*.
- Kaye, K. 2023. The FTC’s ‘profoundly vague’ plan to force companies to destroy algorithms could get very messy. <https://www.protocol.com/enterprise/ftc-algorithm-data-model-ai/>.
- Krizhevsky, A.; Hinton, G.; et al. 2009. Learning multiple layers of features from tiny images.
- Kublik, V. 2023. EU/US Copyright Law and Implications on ML Training Data. <https://valohai.com/blog/copyright-laws-and-machine-learning/>.
- Madry, A.; Makelov, A.; Schmidt, L.; Tsipras, D.; and Vladu, A. 2017. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*.
- Mantelero, A. 2013. The EU Proposal for a General Data Protection Regulation and the roots of the ‘right to be forgotten’. *Computer Law & Security Review*, 29(3): 229–235.
- Nguyen, T. T.; Huynh, T. T.; Nguyen, P. L.; Liew, A. W.-C.; Yin, H.; and Nguyen, Q. V. H. 2022. A survey of machine unlearning. *arXiv preprint arXiv:2209.02299*.
- Paszke, A.; Gross, S.; Massa, F.; Lerer, A.; Bradbury, J.; Chanan, G.; Killeen, T.; Lin, Z.; Gimelshein, N.; Antiga, L.; et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32.
- Shibata, T.; Irie, G.; Ikami, D.; and Mitsuzumi, Y. 2021. Learning with Selective Forgetting. In *IJCAI*, volume 3, 4.
- Simonyan, K.; and Zisserman, A. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*.
- Xiao, H.; Rasul, K.; and Vollgraf, R. 2017. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*.



# Supplementary Materials

## A Datasets

We used the three different datasets as follows:

- **CIFAR-10.** The CIFAR-10 dataset (Krizhevsky, Hinton et al. 2009) is a widely used benchmark in classification tasks. It consists of 60,000 images in ten classes. The dataset is divided into a training set of 50,000 images and a test set of 10,000 images. We experiment to erase only one class (5,000 images) out of 10 classes.
- **Fashion-MNIST.** The Fashion-MNIST dataset (Xiao, Rasul, and Vollgraf 2017) is popular in classification tasks. It contains 70,000 grayscale images of various fashion items, categorized into ten classes. The dataset is divided into a training set of 60,000 images and a test set of 10,000 images. We experiment with erasing only one class (6,000 images) out of 10 classes. We utilize the dataset to evaluate unlearning performance in grayscale images.
- **VGGFace2.** The VGGFace2 dataset (Simonyan and Zisserman 2014) is a large-scale face dataset designed for face recognition tasks. This dataset consists of facial data and is closely related to tasks to preserve privacy. Given the high similarity among classes, it is a crucial dataset for assessing the effectiveness of unlearning methods in real-life scenarios involving facial data. It consists of diverse face images that vary regarding identities, poses, illuminations, backgrounds, and expressions. The dataset contains over 3.31 million images from more than 9,000 individuals. But, to experiment with our unlearning task, we randomly chose ten individuals from a training dataset containing over 600 images, ensuring a balanced distribution of gender.

## B Evaluation Metrics for Unlearning Performance

The results in our experiments are evaluated based on the following metrics:

**Accuracy.** In order to assess a classifier’s performance, accuracy is frequently utilized. It measures the percentage of samples for which the true classes can be predicted with the maximum degree of certainty. Accuracy of a model  $\mathcal{M}_\theta$  tested on a dataset of  $N$  samples  $\{(x_1, y_1), \dots, (x_N, y_N)\}$  is formulated as follows:

$$\text{ACC} = 100 \cdot \frac{\sum_{i=1}^N \delta(\sigma(\mathcal{M}_\theta(x_i)), y_i)}{N}, \quad (9)$$

where  $\delta(\cdot, \cdot)$  is the Kronecker delta function.

**Unlearning Score (US).** Effective unlearning performance refers to the ability of a model to effectively forget information from the forgetting data, while concurrently retaining the relevant information from the retain data. However, determining the most effective metric for unlearning is challenging because of the orthogonal objective between evaluating and measuring forgetting vs. retain data accuracy, where they are not in a linear relationship. For instance,

when comparing two unlearning approaches, one exhibits good performance in forgetting data but does poorly retain data well. At the same time, the other will be the opposite case, with poor performance in forgetting data but great accuracy in retaining data. In such scenarios, it becomes quite challenging to determine which method is better based solely on one of the two accuracies. Therefore, while considering both accuracies is essential, there is no straightforward way to assess both simultaneously. Such discrepancy leads to difficulties in evaluating unlearning performance.

Therefore, to evaluate unlearning performance more accurately and effectively, we propose and define a new metric, called Unlearning Score (US), which effectively characterizes and combines the two accuracies into a single value to assess unlearning performance. Since accuracy is measured in percentages, we normalize it to a range of 0 to 1 by dividing by 100. As accuracy for forgetting data is preferred to be lower, we subtract the value from 1 to convert it into a higher-is-better range. Next, we input the values into the exponential function. The following equation pertains to the retain data:

$$\text{US}_r = \exp\left(\frac{\text{acc}_r}{100}\right), \quad (10)$$

where  $\text{acc}_r$  is accuracy of retain data.

Similarly, the following equation is defined for the forgetting data:

$$\text{US}_f = \exp\left(1 - \frac{\text{acc}_f}{100}\right), \quad (11)$$

where  $\text{acc}_f$  is accuracy of forgetting data.

In fact, we use exponential functions, which offer a good way to assign and map weights to values, much better than linear functions. In other words, rather than simply using the two accuracies as they are, this approach enables us to better characterize higher scores as accuracies increase. Conversely, for lower accuracies, we can assign lower scores. We calculate the average of  $\text{US}_r$  and  $\text{US}_f$  values obtained through Eq. 10 and Eq. 11, respectively. Then, we normalize them to range from 0 to 1 using *min-max* scaling with  $\exp(1)$  and  $\exp(0)$ .

Our final US is constructed and derived to Eq. 8 as follows:

$$\text{US}(\text{acc}_r, \text{acc}_f) = \frac{\text{US}_r + \text{US}_f - \exp(0)}{\exp(1) - \exp(0)} \quad (12)$$

$$= \frac{\text{US}_r + \text{US}_f - 2 \cdot \exp(0)}{2 \cdot (\exp(1) - \exp(0))} \quad (13)$$

$$= \frac{\exp\left(\frac{\text{acc}_r}{100}\right) + \exp\left(1 - \frac{\text{acc}_f}{100}\right) - 2}{2 \cdot (\exp(1) - 1)}$$

By introducing this novel metric, US, we can more effectively characterize and evaluate whether an unlearning method has properly forgotten information from forgetting data, while retained information from retain data simultaneously. Throughout experiments, we show that US effectively characterizes and captures the underlying performance

Table 6: Unlearning performance based on changes in the  $\alpha$  value in knowledge distillation.

$\alpha$		0			0.2			0.5			0.8			1		
Metrics		$D_{tr}$	$D_{tf}$	US	$D_{tr}$	$D_{tf}$	US	$D_{tr}$	$D_{tf}$	US	$D_{tr}$	$D_{tf}$	US	$D_{tr}$	$D_{tf}$	US
CIFAR-10	VGG16	75.31	0	0.8269	91.93	0	0.9386	92.28	0	0.9412	92.17	0	0.9404	92.14	0	0.9402
	ResNet18	92.87	0	0.9456	93.38	0	0.9493	93.50	0	0.9502	92.47	2.4	0.9239	91.58	6.7	0.8849
	ResNet50	91.75	0	0.9374	93.51	0	0.9503	93.43	0	0.9497	90.06	2.8	0.9033	86	10.30	0.8192
	ViT	78.46	0	0.8467	80.65	0	0.8608	81.22	0	0.8645	81.16	0	0.8642	79.11	0.5	0.8469
Fashion-MNIST	VGG16	77.38	0	0.8399	94.21	0	0.9555	93.91	0	0.9532	92.78	0	0.9449	80.47	4.9	0.8218
	ResNet18	91.14	0	0.9329	93.92	0	0.9533	94.6	0	0.9584	93.38	0.6	0.9446	91.54	7.4	0.8794
	ResNet50	85.57	0	0.8937	94.67	0	0.9590	93	0	0.9465	93.18	0.1	0.9471	92.43	8.3	0.8793
	ViT	88.22	0	0.9121	88.38	0	0.9132	88.57	0	0.9146	88.75	0	0.9158	88.44	0	0.9136
VGGFace2	VGG16	91.93	0	0.9386	93.35	0	0.9491	96.04	0	0.9693	96.99	0	0.9765	96.99	0	0.9765
	ResNet18	56.01	0	0.7185	85.12	0	0.8906	94.62	0	0.9585	94.30	0	0.9562	93.19	0	0.9479
	ResNet50	90.82	0	0.9306	93.67	0	0.9514	94.46	0	0.9573	89.39	4.76	0.8836	89.87	14.28	0.8185
	ViT	94.30	0	0.9562	95.56	0	0.9657	95.88	0	0.9681	95.72	0	0.9669	95.56	0	0.9657

Table 7: Unlearning performance based on changes in the  $T$  value in knowledge distillation on CIFAR-10 with ResNet18.

T	1	4	8	16	Original
$D_r$	99.98	99.98	99.97	99.97	99.98
$D_f$	0	0	0	0	100
$D_{tr}$	93.4	93.53	93.32	93.24	93.13
$D_{tf}$	0	0	0	0	96.60
US	0.9495	0.9504	0.9489	0.9483	0.4575

Table 8: Unlearning performance based on changes in the  $x$  value of  $T^x$  in knowledge distillation on CIFAR-10 with ResNet18.

$x$	1	2	3	4	Original
$D_r$	99.97	99.98	99.95	90.63	99.98
$D_f$	0	0	1.94	4.77	100
$D_{tr}$	93.37	93.53	92.67	82.09	93.13
$D_{tf}$	0	0	2.7	5.03	96.60
US	0.9493	0.9504	0.9230	0.8315	0.4575

of forgetting and retain data performance across different proposed methods.

### C Hyper-parameters effects in KD

Table 6 illustrates the variations in unlearning performance based on the hyper-parameter  $\alpha$  in knowledge distillation. When  $\alpha$  is set to 0, our loss function  $\mathcal{L}$  employs only  $\mathcal{L}_{CE}$ , focusing solely on forgetting data. As a result, it may not effectively retain information from the boundary, leading to a potential drop of up to approximately 40%. On the other hand, setting  $\alpha$  to 1 utilizes only  $\mathcal{L}_{DI}$ , prioritizing the retention of the boundary. Although this approach may preserve boundary information well, it might struggle to forget the forgetting data properly. Hence, striking the right balance between forgetting data and boundary information through an appropriate  $\alpha$  value in knowledge distillation is crucial, as shown in Table 6. Table 7 illustrates the variation in unlearning performance based on the hyper-parameter  $T$  in

Table 9: Original PGD vs. Partial-PGD for all datasets.

		Original PGD				Partial-PGD			
Metrics		$D_{tr}$	$D_{tf}$	Time (s)	US	$D_{tr}$	$D_{tf}$	Time (s)	US
CIFAR-10	VGG16	92.03	0	14.18	0.9394	92.18	0	3.76	0.9405
	ResNet18	92.97	0	18.19	0.9463	93.53	0	4.37	0.9504
	ResNet50	91.84	0	44.15	0.9380	93.52	0	7.76	0.9503
	ViT	78.07	0	237.36	0.8442	81.14	0	25.93	0.8640
Fashion-MNIST	VGG16	94.15	0	16.61	0.9551	93.89	0	8.75	0.9531
	ResNet18	94.49	0	21.35	0.9576	94.54	0	5.194	0.9579
	ResNet50	94.47	0	51.74	0.9574	94.48	0	9.14	0.9575
	ViT	87.4	0	23.99	0.9063	87.44	0	13.396	0.9066
VGGFace2	VGG16	96.29	0	19.95	0.9349	96.70	0	5.60	0.9743
	ResNet18	91.42	0	29.21	0.9467	95.34	0	6.51	0.9639
	ResNet50	93.02	0	298.15	0.9712	93.28	0	17.77	0.9485
	ViT	95.76	0	18.65	0.9672	95.5	0	6.748	0.9651

Table 10: Effect of Softmax vs. Double Softmax for all datasets.

		w/o Double Softmax				w/ Double Softmax			
Metrics		$D_{tr}$	$D_{tf}$	Time (s)	US	$D_{tr}$	$D_{tf}$	Time (s)	US
CIFAR-10	VGG16	92.02	0	3.88	0.9472	92.18	0	3.76	0.9405
	ResNet18	93.10	0	4.46	0.9430	93.53	0	4.37	0.9504
	ResNet50	92.53	0	7.56	0.9393	93.52	0	7.76	0.9503
	ViT	78.73	0	69.42	0.8484	81.14	0	25.93	0.8640
Fashion-MNIST	VGG16	84.74	0	10.90	0.8880	93.89	0	8.75	0.9531
	ResNet18	91.42	0.1	25.87	0.9341	94.54	0	5.19	0.9579
	ResNet50	80.91	0	93.49	0.8625	94.48	0	9.13	0.9575
	ViT	87.01	0	61.37	0.9036	87.44	0	13.39	0.9066
VGGFace2	VGG16	92.94	0	3.71	0.9505	96.70	0	5.60	0.9743
	ResNet18	93.54	0	8.75	0.9468	95.34	0	6.51	0.9639
	ResNet50	93.03	0	26.90	0.9461	93.28	0	17.77	0.9485
	ViT	94.91	0	8.49	0.9608	95.50	0	6.74	0.9651

knowledge distillation. In our experiments,  $T = 4$  yielded the best performance; however, variations in  $T$  showed a difference in accuracy of 0.2% as indicated in Table 7 under  $D_{tr}$ . Table 8 illustrates the variation in unlearning performance based on the exponent of  $T^x$  in Eq. 7,  $T$  is fixed at 4. In our experiments,  $x = 2$  yielded the best performance, whereas values greater than 3 demonstrated poorer performance. Experiments with hyper-parameters tuning show that appropriately selecting values in knowledge distillation can yield the better performance in the unlearning task.

### D Original PGD vs. Partial-PGD

We conduct experiments on various models and datasets to demonstrate the temporal efficiency and performance advantage of Partial-PGD. In Table 9, the original PGD also presents an excellent performance in terms of unlearning.

However, we show that Partial-PGD exhibits comparable or superior performance to original PGD, notably in VG-GFace2 with ResNet50. On the other hand, it can save unlearning process time up to 16.77 times. The original PGD requires more time as the model has to utilize the complete model layers. In contrast, as depicted in Fig. 1, Partial-PGD can be considered more effective, as it only uses particular layers to achieve the desired objectives faster.

## E Effectiveness of Double Softmax

As shown in Eq. 5, double Softmax provides performance robustness across various datasets and models. In Table 10, we conduct experiments to examine the effects of double Softmax across different datasets and models. Overall, double Softmax facilitates a faster unlearning convergence speed. Furthermore, though the difference is marginal, our experimental results demonstrate higher accuracy performance across most models. Especially in the case of Fashion-MNIST, notable improvements can be observed. Double Softmax generates softer logits, enhancing robustness against outliers of adversarial examples and improving training stability.

## F Additional Ablation on Data Usage Ratio

We conduct an Ablation Study to investigate whether reducing the amount of randomly selected forgetting data involved in our algorithm’s unlearning process impacts performance while maintaining the possibility of unlearning. Table 11 presents the results when reducing the data used in the unlearning process across various datasets. The remarkable finding is that even with a reduction in the quantity of forgetting data, there is no significant decline in performance from an accuracy perspective. Additionally, a decrease in the completion time of the unlearning process can also be observed. It can be observed that for ViT on Fashion-MNIST, the accuracy of  $D_{tf}$  remains at 0.1%.

## G Unlearning Performance on Every Class

We conduct experiments on our method for all classes to showcase its robust performance regardless of datasets and classes. Table 12 presents experiments on CIFAR-10, demonstrating our method’s ability to quickly erase an entire class, while retaining other information in as little as 2.5 seconds. Table 13 shows experiments on Fashion-MNIST, where although perfect erasure of a single class might not always be achieved, our method consistently demonstrates efficient and effective performance across all other experiments. Finally, Table 14 highlights experiments on VG-GFace2, showing our method’s remarkable performance even on face datasets with high inter-class similarity.

Table 11: Unlearning performance with varying amounts of data used for unlearning.

Model		VGG16			ResNet18			ResNet50			ViT		
Total Extra Data Used		100%	50%	10%	100%	50%	10%	100%	50%	10%	100%	50%	10%
CIFAR-10	$D_{tr}$	92.18	92.42	92.38	93.53	93.51	93.38	93.52	93.63	93.37	81.14	81.14	81.60
	$D_{tf}$	0	0	0	0	0	0	0	0	0	0	0	0
	Time	3.76	1.91	1.21	4.37	2.28	1.45	7.76	3.81	1.62	25.93	25.63	14.55
	US	0.9405	0.9422	0.9420	0.9504	0.9503	0.9493	0.9503	0.9512	0.9493	0.8640	0.8640	0.8662
Fashion-MNIST	$D_{tr}$	93.89	94.23	93.74	94.54	94.67	97.19	94.48	94.21	84.88	87.44	87.09	87.46
	$D_{tf}$	0	0	0	0	0	0	0	0	0	0	0	0.1
	Time	8.75	2.20	0.48	5.19	1.96	0.63	9.14	4.54	1.04	13.39	4.88	2.69
	US	0.9531	0.9556	0.9520	0.9579	0.9589	0.9487	0.9575	0.9555	0.8890	0.9066	0.9042	0.9060
VGGFace2	$D_{tr}$	96.70	95.83	95.88	95.34	94.35	94.46	93.28	94.24	93.13	95.50	95.82	95.88
	$D_{tf}$	0	0	0	0	0	0	0	0	0	0	0	0
	Time	5.60	5.12	5.36	6.51	4.22	1.8	17.77	23.09	15.35	6.74	2.46	2.04
	US	0.9743	0.9677	0.9681	0.9639	0.9565	0.9573	0.9485	0.9557	0.9475	0.9651	0.9676	0.9680

Table 12: Unlearning performance for each class on CIFAR-10

Forgetting Class		0	1	2	3	4	5	6	7	8	9
VGG16	$D_r$	99.98	99.98	99.98	99.98	99.98	99.98	99.98	99.98	99.97	99.97
	$D_f$	0	0	0	0	0	0	0	0	0	0
	$D_{tr}$	92.7	92.2	93.35	94.44	93.17	93.98	93.54	92.47	92.24	92.35
	$D_{tf}$	0	0	0	0	0	0	0	0	0	0
	Time (s)	3.75	7.31	3.7	3.71	3.68	3.66	2.5	3.68	3.65	3.72
	US	0.9443	0.9406	0.9491	0.9572	0.9477	0.9537	0.9431	0.9426	0.9409	0.9417
ResNet18	$D_r$	99.97	99.98	99.98	99.98	99.97	99.97	99.98	99.98	99.98	99.98
	$D_f$	0	0	0	0	0	0	0	0	0	0
	$D_{tr}$	93.84	93.44	94.39	95.26	93.86	94.53	93.53	93.48	93.58	93.51
	$D_{tf}$	0	0	0	0	0	0	0	0	0	0
	Time (s)	4.37	4.42	4.44	4.41	3.52	3.25	3.32	4.40	4.42	4.47
	US	0.9527	0.9497	0.9568	0.9633	0.9528	0.9578	0.9504	0.9500	0.9508	0.9502
ResNet50	$D_r$	99.94	99.93	99.94	99.94	99.94	99.97	99.94	99.92	99.93	99.89
	$D_f$	0	0	0	0	0	0	0	0	0	0
	$D_{tr}$	93.93	93.07	94.45	95.26	94.06	94.54	93.56	93.48	93.58	92.97
	$D_{tf}$	0	0	0	0	0	0	0	0	0	0
	Time (s)	7.42	7.41	7.52	7.60	7.49	7.86	7.49	7.54	7.47	7.47
	US	0.9534	0.9470	0.9572	0.9633	0.9543	0.9579	0.9506	0.9500	0.9508	0.9463
ViT	$D_r$	88.43	88.2	88.42	89.74	87.83	88.96	86.48	86.72	87.52	88.07
	$D_f$	0	0	0.02	0	0	0	0	0	0	0
	$D_{tr}$	82.68	81.60	83.31	84.10	82.14	82.65	80.73	80.84	81.13	82.07
	$D_{tf}$	0	0	0	0	0	0	0	0	0	0
	Time (s)	15.20	16.68	250.82	21.16	46.01	33.49	104.40	33.68	25.20	8.65
	US	0.8742	0.8670	0.8776	0.8837	0.8706	0.8732	0.8613	0.8620	0.8639	0.8701

Table 13: Unlearning performance for each class on Fashion-MNIST

Forgetting Class		0	1	2	3	4	5	6	7	8	9
VGG16	$D_r$	99.84	99.75	99.74	99.62	99.84	99.2	99.88	99.88	99.78	99.85
	$D_f$	0	0	0	0	0	0	0	0	0	0
	$D_{tr}$	96.02	94.1	95.27	94.93	95.37	93.53	97.3	94.83	94.31	94.5
	$D_{tf}$	0	0	0	0	0	0	0	0	0	0
	Time (s)	4.35	8.56	4.24	4.32	4.3	4.36	4.33	4.36	4.35	4.33
	US	0.9691	0.9546	0.9634	0.9608	0.96422	0.9504	0.9789	0.9601	0.9562	0.9576
ResNet18	$D_r$	97.44	97.21	97.94	97.65	96.35	97.05	98.72	98.31	97.9	98.41
	$D_f$	0	0	0	0	0	0	0	0	0	0
	$D_{tr}$	94.94	93.75	95.36	94.63	93.84	93.48	96.71	94.94	94.48	95.08
	$D_{tf}$	0	0	0	0	0	0	0	0	0	0
	Time (s)	7.77	15.87	53.49	25.62	15.44	20.89	15.32	10.24	5.15	5.22
	US	0.9609	0.9520	0.96415	0.9578	0.9527	0.9500	0.9743	0.9609	0.9575	0.9620
ResNet50	$D_r$	97.47	98.02	96.45	98.05	97.22	98.16	98.35	98.13	98.20	98.27
	$D_f$	0.05	0	0	0	0	0	0	0	0	0
	$D_{tr}$	94.94	94.26	94.11	95.23	94.76	94.53	96.37	94.62	94.63	95.05
	$D_{tf}$	0	0	0	0.1	0	0	0	0	0	0
	Time (s)	105.24	118.87	99.68	9.14	99.77	8.83	25.99	17.50	10.86	9.33
	US	0.9609	0.9558	0.9547	0.9631	0.9596	0.9578	0.9718	0.9585	0.9586	0.9617
ViT	$D_r$	93.03	90.52	91.02	91.72	92.82	89.82	93.22	91.99	89.82	91.63
	$D_f$	0	0	0	0	0	0	0	0	0	0
	$D_{tr}$	90.37	87.63	88.72	89.23	90.74	86.91	91.45	88.93	87.17	88.54
	$D_{tf}$	0	0	0	0	0	0	0	0	0	0
	Time (s)	9.88	4.93	19.8	9.77	19.7	4.95	14.71	4.93	9.82	4.91
	US	0.9273	0.9079	0.9156	0.9192	0.9300	0.9029	0.9351	0.9171	0.9047	0.9143

Table 14: Unlearning performance for each class on VGGFace2

Forgetting Class		0	1	2	3	4	5	6	7	8	9
VGG16	$D_r$	99.71	99.89	99.78	99.59	99.52	99.38	99.87	99.74	99.44	99.73
	$D_f$	0	0	0	0	0	0	0	0	0	0
	$D_{tr}$	96.01	97.11	96.39	96.34	95.87	96.03	97.12	96.25	95.72	96.82
	$D_{tf}$	0	0	0	0	0	0	0	0	0	0
	Time (s)	19.27	20.01	15.25	5.92	6.35	5.74	6.45	7.63	5.86	6.42
	US	0.9690	0.9774	0.9719	0.9715	0.9679	0.9692	0.9775	0.9708	0.9668	0.9752
ResNet18	$D_r$	99.85	99.85	99.60	99.82	99.68	99.79	99.54	99.83	99.72	99.94
	$D_f$	0	0	0	0	0	0	0	0	0	0
	$D_{tr}$	94.58	94.70	94.59	95.23	95.39	95.08	95.21	95.11	95.25	95.55
	$D_{tf}$	0	0	0	0	0	0	0	0	0	0
	Time (s)	17.94	19.22	10.91	16.62	8.81	8.27	9.1	10.64	8.38	8.84
	US	0.9582	0.9591	0.958	0.9631	0.9643	0.9620	0.9630	0.9622	0.9633	0.9655
ResNet50	$D_r$	95.19	98.51	98.42	97.84	98.07	98.67	98.62	98.61	98.72	98.61
	$D_f$	0.05	0	0	0	0	0	0	0	0	0
	$D_{tr}$	95.88	93.9	93.94	93.64	93.49	94.61	94.73	94.46	94.62	94.76
	$D_{tf}$	0	0	0	0	0	0	0	0	0	0
	Time (s)	22.78	20.06	22.17	17.57	17.65	16.84	18.26	21.6	16.92	18.74
	US	0.9680	0.9531	0.9534	0.9512	0.9501	0.9584	0.9593	0.9573	0.9585	0.9596
ViT	$D_r$	94.88	94.93	95.23	94.77	95.09	94.92	95.91	95.27	95.33	95.22
	$D_f$	0	0	0	0	0	0	0	0	0	0
	$D_{tr}$	95.38	95.02	94.76	95.23	95.55	95.40	95.21	95.43	95.88	95.23
	$D_{tf}$	0	0	0	0	0	0	0	0	0	0
	Time (s)	6.19	6.23	6.93	4.94	5.01	9.99	6.02	6.03	5.31	5.72
	US	0.9642	0.9615	0.9596	0.9631	0.9655	0.9644	0.9630	0.9646	0.96802	0.9631