

BLIND UNLEARNING: UNLEARNING WITHOUT A FORGET SET

Anonymous authors

Paper under double-blind review

ABSTRACT

Machine unlearning is the study of methods to efficiently remove the influence of some subset of the training data from the parameters of a previously-trained model. Existing methods typically require direct access to the “forget set” – the subset of training data to be forgotten by the model. This limitation impedes privacy, as organizations need to retain user data for the sake of unlearning when a request for deletion is made, rather than being able to delete it immediately. We first introduce the setting of *blind unlearning* – unlearning without explicit access to the forget set. Then, we propose a method for approximate unlearning called RELOAD, that leverages ideas from gradient-based unlearning and neural network sparsity to achieve blind unlearning. The method serially applies an ascent step with targeted parameter re-initialization and fine-tuning, and on empirical unlearning tasks, RELOAD often approximates the behaviour of a from-scratch retrained model better than approaches that leverage the forget set. Finally, we extend the blind unlearning setting to *blind remedial learning*, the task of efficiently updating a previously-trained model to an amended dataset¹.

1 INTRODUCTION

Machine unlearning poses the problem of removing the influence of certain instances in the training data on a given statistical model (Bourtoule et al., 2019). Motivated by “right to be forgotten” provisions, like those in the European Union’s General Data Protection Regulation (GDPR) (European Parliament & Council of the European Union), methods in machine unlearning aim to provide efficient means to “forget” specific data points from a trained model without requiring that it be retrained from scratch. As larger models become more prevalent (Achiam et al., 2023), the need to unlearn specific data instances without retraining from scratch is increasingly important.

Contemporary unlearning methods generally require explicit access to the so-called “forget set” – the subset of training data to be forgotten by the model. For example, one approach entails performing steps of gradient ascent on the loss landscape characterized by the forget set in order to remove its influence on the model weights (Thudi et al., 2022). However, the reliance of these methods on the forget set introduces a tension in the context of preserving user privacy: in order to enable unlearning, organizations must retain the complete original set of user data on which the model was trained. The retention of this data, even for the purpose of unlearning, can expose organizations and individuals to risks associated with data breaches or unauthorized access. To bridge this gap, there is a clear need for unlearning methods that operate without requiring access to the forget set. Existing work aims to reduce the reliance on the forget set, but is limited to the constrained task of forgetting classes of data, and requires knowing which classes are being forgotten (Tarun et al., 2023).

This work presents an algorithm for machine unlearning in the absence of an explicitly defined forget set; a setting we establish as “blind unlearning.” Our method, RELOAD, assumes that the modeller only has access to (a) a model trained on a dataset \mathcal{D} , (b) the “retain set,” $\mathcal{D}_{new} \triangleq \mathcal{D} \setminus \mathcal{D}_{forget}$, and (c) cached gradients from the last iteration of training on \mathcal{D} . Notable in its absence from these requirements is the forget set – this means that RELOAD allows for deletion of instances in \mathcal{D}_{forget} at the conclusion of training without inhibiting downstream unlearning.

In this vein, our work makes the following contributions:

¹A software implementation of our work can be found in [this code repository](#).

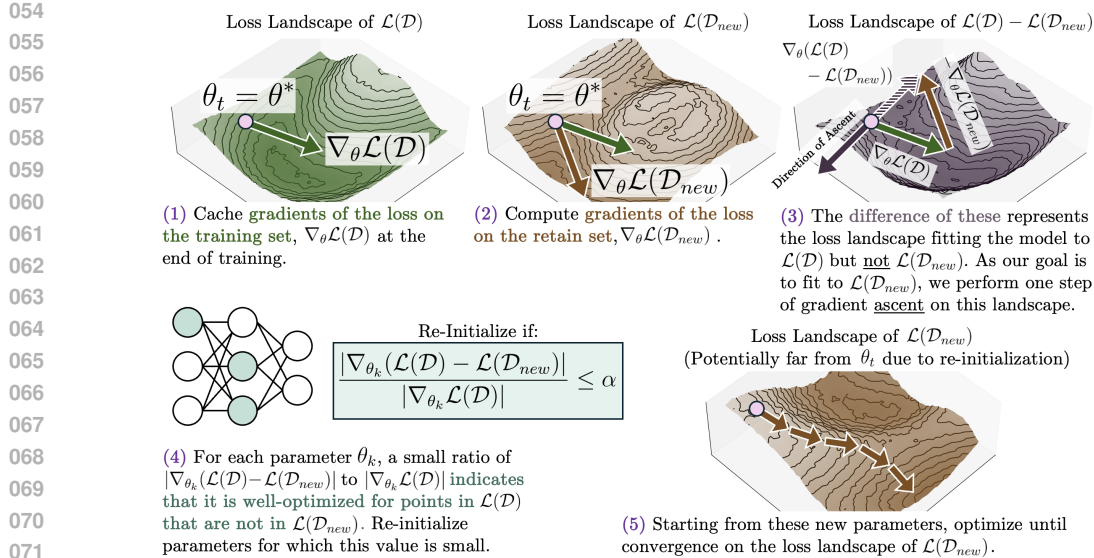


Figure 1: High-level overview of the RELOAD algorithm for blind approximate unlearning and remedial learning. The algorithm marries a gradient-based unlearning step modified for the blind unlearning setting (Steps (1) through (3)) with a weight saliency-based selective re-initialization (Step (4)) and subsequent fine-tuning (Step (5)). Because the blind unlearning setting prohibits taking gradients with respect to \mathcal{D}_{forget} , RELOAD exploits the linearity of differentiation to treat $\nabla_{\theta}(\mathcal{L}(\mathcal{D}) - \mathcal{L}(\mathcal{D}_{new}))$ as a proxy for $\nabla_{\theta}\mathcal{L}(\mathcal{D}_{forget})$ at the location in parameter space corresponding to θ_t . This allows us to apply one gradient step in this direction. Intuitively, this update in Step (3) removes information about \mathcal{D}_{forget} from *all* network parameters, while the re-initialization in Step (4) re-initialises those parameters with a uniquely strong correspondence to \mathcal{D}_{forget} (for which a single ascent step will not fully remove this information). RELOAD achieves state-of-the-art performance on a collection of unlearning tasks, often outperforming baselines with direct access to \mathcal{D}_{forget} .

1. We introduce RELOAD, an algorithm for approximate blind unlearning in parametric models. RELOAD marries ideas from gradient-based unlearning algorithms and neural network sparsity to achieve blind unlearning. We formally show that the requirements of the RELOAD algorithm satisfy the blind unlearning by not permitting recoverability of instances in the forget set in the common setting of softmax classification.
2. Empirical evaluations demonstrate that RELOAD effectively tackles machine unlearning in several diverse settings, often faithfully approximating the behaviour of a from-scratch retrained model *better than existing unlearning approaches that leverage the forget set*.
3. We extend the RELOAD framework to the setting of “remedial learning,” which aims to efficiently update a statistical model given that some instances in the training data that have been amended since the model was originally trained. This enables computationally-efficient data correction in pretrained models without the need for costly retraining.

2 BACKGROUND

2.1 SETTING AND NOTATION

Let $\mathcal{D} = \{(X_i, Y_i)\}_{i=1, \dots, N}$ represent a collection of i.i.d. data, where $X \in \mathcal{X}$ represents input covariates and $Y \in \mathcal{Y}$ represents labels for supervised learning. Then, for some class of models \mathcal{M} , let θ^* represent the parameters that minimize the empirical loss with respect to training data \mathcal{D} ,

$$\theta^* \triangleq \arg \min_{\theta \in \Theta} \mathbb{E}_{(X_i, Y_i) \sim \mathcal{D}} \mathcal{L}((X_i, Y_i); \theta). \quad (1)$$

We denote an instantiation of \mathcal{M} trained on \mathcal{D} as $\mathcal{M}^{(\theta^*)}$. After $\mathcal{M}^{(\theta^*)}$ is trained, assume that some transformation is applied to \mathcal{D} to yield \mathcal{D}_{new} (e.g., deleting instances from \mathcal{D} that are in \mathcal{D}_{forget}). Then, θ^\sim represents the parameters that minimize the empirical loss with respect to \mathcal{D}_{new} ,

$$\theta^\sim \triangleq \arg \min_{\theta \in \Theta} \mathbb{E}_{(X_i, Y_i) \sim \mathcal{D}_{new}} \mathcal{L}((X_i, Y_i); \theta). \quad (2)$$

Our general goal, encompassing both unlearning and remedial learning, is transforming $\mathcal{M}^{(\theta^*)}$ into $\mathcal{M}^{(\theta^\sim)}$ without naively obtaining $\mathcal{M}^{(\theta^\sim)}$ by re-training an instance of \mathcal{M} on \mathcal{D}_{new} from scratch.

Machine Unlearning. In the machine unlearning setting, the transformation of \mathcal{D} into \mathcal{D}_{new} consists of first identifying a subset of the data whose influence to remove, \mathcal{D}_{forget} , and taking $\mathcal{D}_{new} \triangleq \mathcal{D} \setminus \mathcal{D}_{forget}$. These remaining instances represent the portion of the training data that is retained – the full training set, less those instances to be forgotten. The goal of approximate unlearning methods – of which RELOAD is one (see Section 2.2) – is to efficiently learn an approximation of $\mathcal{M}^{(\theta^\sim)}$. The classical setting assumes that the modeller has access to the trained model $\mathcal{M}^{(\theta^*)}$, the training dataset \mathcal{D} , the remaining data \mathcal{D}_{new} , and the forget set \mathcal{D}_{forget} (Cao & Yang, 2015).

Remedial Learning. The unlearning setting is subject to the restriction that $\mathcal{D}_{new} \triangleq \mathcal{D} \setminus \mathcal{D}_{forget}$; however, one may also consider the broader setting wherein \mathcal{D}_{new} is the result of some arbitrary item-wise transformation to \mathcal{D} . Formally, let $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{X} \times \mathcal{Y}$ denote a transformation, and write $(X'_i, Y'_i) = f(X_i, Y_i)$. Then, \mathcal{D}_{new} represents the result of applying f item-wise to K elements of \mathcal{D} , and applying the identity transform to the remaining $N - K$ elements, as

$$\mathcal{D}_{new} = \{(X'_i, Y'_i)\}_{i=1, \dots, K} \cup \{(X_i, Y_i)\}_{i=K+1, \dots, N}. \quad (3)$$

This represents a generalization of the unlearning problem, as we wish to “unlearn” the influence of $\{(X_i, Y_i)\}_{i=1, \dots, K}$ on our original model, and “relearn” the influence of $\{(X'_i, Y'_i)\}_{i=1, \dots, K}$. This setting encompasses the following data transformations, among others:

1. *Covariate Correction:* $\mathcal{D}_{new} = \{(X'_i, Y_i)\}_{i=1, \dots, K} \cup \{(X_i, Y_i)\}_{i=K+1, \dots, N}$, where X'_i represents a corrected version of the features X_i , and indices $K + 1, \dots, N$ correspond to those with erroneous covariates (e.g., data was corrupted during collection/pre-processing).
2. *Label Correction:* $\mathcal{D}_{new} = \{(X_i, Y'_i)\}_{i=1, \dots, K} \cup \{(X_i, Y_i)\}_{i=K+1, \dots, N}$, where Y'_i represents a corrected version of the label Y_i , and indices $K + 1, \dots, N$ correspond to those that were originally mis-labelled during annotation.
3. *Backdoor Removal:* $\mathcal{D}_{new} = \{(X'_i, Y_i)\}_{i=1, \dots, K} \cup \{(X_i, Y_i)\}_{i=K+1, \dots, N}$, where X'_i represents a version of the features X_i lacking the injected backdoor pattern, and indices $K + 1, \dots, N$ correspond to those that were originally transformed with a backdoor during processing. Models trained with backdoors in the training set learn shortcuts (Geirhos et al., 2020), which can be exploited to induce misclassification.

This work studies how the RELOAD algorithm accomplishes tasks both in the unlearning setting, and in the setting of remedial learning.

Blind Unlearning / Blind Remedial Learning. In contrast to the classical unlearning (and remedial learning) settings, in which the modeller has access to the forget set, our setting assumes no such access. We call this setting *blind unlearning* (or *blind remedial learning*). The blind unlearning / remedial learning setting has access to the trained model $\mathcal{M}_{\mathcal{D}}$, the new dataset \mathcal{D}_{new} , and (potentially) some limited information about the original data, $\mathcal{I}_{\mathcal{D}}$, from which neither \mathcal{D}_{forget} (in the unlearning setting) or \mathcal{D} (in the remedial learning setting) can be fully reconstructed.

2.2 RELATED WORK

Exact and Approximate Unlearning. Exact unlearning refers to the subclass of algorithms that provide formal guarantees of the extent to which information about \mathcal{D}_{forget} was removed from the weights of a model. The trivial method for exact unlearning consists of naively retraining the model from scratch (this is considered the gold-standard for machine unlearning; see Cao & Yang (2015); Thudi et al. (2022); Shaik et al. (2024)). Other exact unlearning methods include SISA (Bourtoule et al., 2019), which partitions the data to accelerate retraining, Certified Data Removal

(Guo et al., 2019), which performs a Newton update in the opposite direction of the gradient with respect to \mathcal{D}_{forget} , and Certified Graph Unlearning (Chien et al., 2022), which builds on Certified Data Removal using the graph topology to enforce guaranteed unlearning. Unlike exact unlearning methods, approximate unlearning algorithms (like RELOAD) aim to recover the behaviour of a model naively retrained on the new set without providing any formal theoretical guarantees. These methods can be sub-classified into either gradient-based or weight-saliency based algorithms.

Gradient-Based Approximate Unlearning. Many existing approximate unlearning algorithms perform an optimization procedure on $\mathcal{M}(\theta^*)$ using the forget set \mathcal{D}_{forget} and the retain set \mathcal{D}_{new} . One simple standard approach applies gradient ascent on the loss with respect to \mathcal{D}_{forget} , in order to undo the parameter updates induced by those instances during training (Graves et al., 2021; Thudi et al., 2022). Another gradient-based approach leverages a teacher-student method: “Bad Teacher” performs knowledge distillation based on one trained model on \mathcal{D}_{new} (the “good teacher”) and one a randomly initialised model on \mathcal{D}_{forget} (the “bad teacher”) (Chundawat et al., 2022); SCalable Remembering and Unlearning unBound (SCRUB) similarly distills a student model from a teacher trained on \mathcal{D} , but the student learns to selectively disobey the teacher by directly maximizing the loss on \mathcal{D}_{forget} (Kurmanji et al., 2023). A third family directly manipulates the structure of the learned representation space using gradients: Distance-based Unlearning via Centroid Kinematics (DUCK) (Cotogni et al., 2023) drives representations of elements in \mathcal{D}_{forget} towards the nearest incorrect centroid in the feature space, while Boundary Unlearning (Chen et al., 2023) implements class-level unlearning by shifting the decision boundary corresponding to the class(es) defining \mathcal{D}_{forget} .

Weight Saliency-Based Approximate Unlearning. Another class of approximate unlearning methods derives from the hypothesis that identifiable substructures in neural networks often correspond to different subsets of the training data (Pfeiffer et al., 2023). These methods leverage ideas from neural sparsity (Frankle & Carbin, 2018; Chen et al., 2024) to perform targeted unlearning on specific parameters. Saliency unlearning (SalUn) uses a threshold on $\nabla_{\theta} \mathcal{L}(\mathcal{D}_{forget})$ to identify parameters containing the most signal about \mathcal{D}_{forget} and focuses model updates on these parameters (Fan et al., 2023). Selective Synaptic Dampening (SSD) (Foster et al., 2023) extends this idea to avoid gradient steps by scaling parameters based on their Fisher Information Matrix importance scores.

Blind Unlearning. This setting involves unlearning without access to \mathcal{D}_{forget} at the instanced of unlearning. It then reduces to taking a model fit on one dataset \mathcal{D} and adapting it to fitting a new dataset \mathcal{D}_{new} . This connects to domain adaptation, in which differences in datasets may not be explicitly defined. In blind unlearning it is not available. An unlearning baseline, Finetuning (FT) (Warnecke et al., 2023) on the retain-set \mathcal{D}_{new} fulfills the blind criteria. Catastrophically forgetting the last k layers (CF- k) and Exact-unlearning the last k layers (EU- k) (Goel et al., 2022) are also blind. Fisher Forgetting (Fisher) (Golatkar et al., 2020) is also a blind unlearning algorithm, but is theoretically bound by class unlearning. Both FT and CF- k provide no strong theoretical indication of unlearning. EU- k does by re-initialising the last k layers of the model. Our method, RELOAD, provides a stronger indication by selectively re-initialises parameters which know the most about the knowledge we wish to remove.

3 RELOAD

3.1 ALGORITHM REQUIREMENTS

Assumption 1 (Unlearning from Cached Gradients). *In RELOAD, $\mathcal{I}_{\mathcal{D}} = \nabla_{\theta} \mathcal{L}(\mathcal{D})$.*

The following lemma demonstrates why this is a valid choice of $\mathcal{I}_{\mathcal{D}}$ for blind unlearning / remedial learning in the softmax classification setting, because this choice does not permit recovery of the removed instances within \mathcal{D}_{forget} (or of instances in \mathcal{D} and not in \mathcal{D}_{new} in remedial learning) in the common setting of softmax classification.

Definition 1 (Recoverability). *Consider some data, $\mathcal{D} \in \mathcal{D}$, and consider a transformation $f : \mathcal{D} \rightarrow \mathcal{Q}$ that maps \mathcal{D} into an arbitrary output space \mathcal{Q} . \mathcal{D} is recoverable if f is injective.*

Lemma 1 (\mathcal{D}_{forget} is Not Recoverable from $\mathcal{I}_{\mathcal{D}}$ in Softmax Classification). *Consider softmax classification over C classes, where each $Y_i \in [0, 1]^C$ represents a one-hot encoded vector of class labels, $\hat{Y}_i = [e^{z_{i1}} / \sum_{j=1}^C e^{z_{j1}}, \dots, e^{z_{iC}} / \sum_{j=1}^C e^{z_{jC}}]$ represents predicted probabilities for each class*

generated from model logits $Z_i \in \mathbb{R}^C$, and $\mathcal{L}((X_i, Y_i), \hat{Y}_i) = -\sum_{i=1}^N \sum_{j=1}^C Y_{ij} \log \hat{Y}_{ij}$. The transformation $\mathcal{G} : \mathcal{D} \rightarrow \Theta$ s.t. $\mathcal{G}(\mathcal{D}) = \nabla_{\theta} \mathcal{L}(\mathcal{D}) \triangleq \mathcal{I}_{\mathcal{D}}$ is not injective.

Proof. Recall from Section 3.2 that we can write $\nabla_{\theta} \mathcal{L}(\mathcal{D}) - \nabla_{\theta} \mathcal{L}(\mathcal{D}_{new})$ as $\nabla_{\theta} \mathcal{L}(\mathcal{D}_{forget})$. Then, if $|\mathcal{D}_{forget}| > 1$, the numerator $\nabla_{\theta_k} \mathcal{L}(\mathcal{D}_{forget})$ can be written as $\sum_{(X_i, Y_i) \in \mathcal{D}_{forget}} \nabla_{\theta_k} \mathcal{L}((X_i, Y_i), \hat{Y}_i)$. Because summation is not injective, \mathcal{G} is also non-injective. In the other case, if $|\mathcal{D}_{forget}| = 1$, we write $\mathcal{D}_{forget} = \{(X_1, Y_1)\}$. Without loss of generality, $Y_{1j} = 1$ and $Y_{1k} = 0$ for all $k \neq j$. Given $\nabla_{\theta_k} \mathcal{L}((X_1, Y_1), \hat{Y}_1) = -\nabla_{\theta_k} \sum_{i=1}^C Y_{1i} \log \hat{Y}_{1i} = \nabla_{\theta_k} \log \hat{Y}_{1j} = \frac{1}{\hat{Y}_{1j}}$, we can recover the j 'th output of the model, $\hat{Y}_{1j} \cdot \hat{Y}_{1j} = \frac{e^{Z_{1j}}}{\sum_{i=1}^C e^{Z_{1i}}}$. For any element Z_{1k} of Z_1 , $e^{Z_{1k}} = \hat{Y}_{1j} \cdot \sum_{i=1}^C e^{Z_{1i}}$, this implies that $Z_{1k} = \log(\hat{Y}_{1j} \cdot \sum_{i=1}^C e^{Z_{1i}}) = \log(\hat{Y}_{1j}) + \log(\sum_{i=1}^C e^{Z_{1i}})$ which cannot be calculated without knowing the other elements of Z_1 . Thus, given only \hat{Y}_{1j} , no elements of Z_1 can be obtained, hence \mathcal{G} is also injective in this case. \square

3.2 ALGORITHM INTUITION

Direction of Movement. The central challenge of blind unlearning is that taking repeated gradients of $\mathcal{L}(\mathcal{D}_{forget})$ is impossible without access to \mathcal{D}_{forget} . However, from cached gradients of \mathcal{D} at the conclusion of model training, $\nabla_{\theta} \mathcal{L}(\mathcal{D})$, we can infer $\nabla_{\theta} \mathcal{L}(\mathcal{D}_{forget})$.

To do so, let $\hat{Y}_i = \mathcal{M}^{(\theta)}(X_i)$ represent the model's prediction. Then,

$$\nabla_{\theta} \mathcal{L}(\mathcal{D}_{forget}) = \sum_{(X_i, Y_i) \in \mathcal{D}_{forget}} \nabla_{\theta} \mathcal{L}((X_i, Y_i), \hat{Y}_i) = \sum_{(X_i, Y_i) \in \mathcal{D} \setminus \mathcal{D}_{new}} \nabla_{\theta} \mathcal{L}((X_i, Y_i), \hat{Y}_i) \quad (4)$$

where the second equality follows from $\mathcal{D}_{new} = \mathcal{D} \setminus \mathcal{D}_{forget}$. Equivalently,

$$= \sum_{(X_i, Y_i) \in \mathcal{D}} \nabla_{\theta} \mathcal{L}((X_i, Y_i), \hat{Y}_i) - \mathbb{1}_{(X_i, Y_i) \in \mathcal{D}_{new}} \left[\nabla_{\theta} \mathcal{L}((X_i, Y_i), \hat{Y}_i) \right] \quad (5)$$

$$= \sum_{(X_i, Y_i) \in \mathcal{D}} \nabla_{\theta} \mathcal{L}((X_i, Y_i), \hat{Y}_i) - \sum_{(X_i, Y_i) \in \mathcal{D}_{new}} \nabla_{\theta} \mathcal{L}((X_i, Y_i), \hat{Y}_i) \quad (6)$$

$$= \nabla_{\theta} \mathcal{L}(\mathcal{D}) - \nabla_{\theta} \mathcal{L}(\mathcal{D}_{new}). \quad (7)$$

Therefore, a gradient-based descent update in the direction of $\nabla_{\theta} \mathcal{L}(\mathcal{D}_{forget})$ moves the model parameters such that they better fit to \mathcal{D}_{forget} ; because our goal is *unlearning* \mathcal{D}_{forget} , RELOAD instead begins with a **single gradient ascent update step** in this direction.

In unlearning, our goal is to obtain a gradient in the direction of \mathcal{D}_{forget} . The remedial learning case is more general: the goal is to obtain $\nabla_{\theta} \mathcal{L}(\mathcal{D} \cap \mathcal{D}_{new}^c)$, a gradient pointing towards the empirical minimum of the loss on elements that are uniquely contained in \mathcal{D} and not in \mathcal{D}_{new} , and $-\nabla_{\theta} \mathcal{L}(\mathcal{D}^c \cap \mathcal{D}_{new})$, a gradient pointing away from the empirical minimum of the loss on elements uniquely contained in \mathcal{D}_{new} but not in \mathcal{D} . Unlearning represents the special case of this framework in which $\mathcal{D} \cap \mathcal{D}_{new}^c = \mathcal{D}_{forget}$ and $\mathcal{D}^c \cap \mathcal{D}_{new} = \emptyset$. In the remedial learning setting, the desired gradient is also $\nabla_{\theta} \mathcal{L}(\mathcal{D}) - \nabla_{\theta} \mathcal{L}(\mathcal{D}_{new})$; the derivation can be found in Appendix A.1. This informs Step (2-3) in Figure 1.

Targeted Parameter Adjustments. Taking a gradient step in this direction, however, is insufficient for unlearning (or remedial learning) for two reasons. First, we are limited to a single gradient step in this direction (Assumption 1), and second, theory from network modularity (Rodriguez et al., 2019) suggests that a small subset of parameters contain a disproportionate amount of the necessary information to characterize instances in \mathcal{D}_{forget} . While one ascent step may be useful at removing what little information about \mathcal{D}_{forget} is included across all network parameters, it is less plausible that this single step will remove information about \mathcal{D}_{forget} from the subset of parameters most responsible for its characterization.

We therefore perform selective re-initialization of these parameters as follows. Consider the gradient $\nabla_{\theta_k} \mathcal{L}(\mathcal{D}_{forget})$, the gradient of the loss with respect to instances in \mathcal{D}_{forget} and with respect to a

particular parameter θ_k . If this gradient is small, it means that θ_k is well-optimized to characterize instances in \mathcal{D}_{forget} ; if this gradient is large, it means that θ_k poorly characterizes these instances. Although the absolute values of these gradients are largely meaningless, the relative magnitude of $\nabla_{\theta_k} \mathcal{L}(\mathcal{D}_{forget})$ compared $\nabla_{\theta_k} \mathcal{L}(\mathcal{D})$ is a meaningful representation of the extent to which θ_k is responsible for characterizing information about \mathcal{D}_{forget} . We call this the *knowledge value* of parameter θ_k , and formally define it as,

$$KV_{\theta_k} \triangleq \frac{|\nabla_{\theta_k} \mathcal{L}(\mathcal{D}_{forget})| + \epsilon}{|\nabla_{\theta_k} \mathcal{L}(\mathcal{D})| + \epsilon} = \frac{|\nabla_{\theta_k} \mathcal{L}(\mathcal{D}) - \nabla_{\theta_k} \mathcal{L}(\mathcal{D}_{new})| + \epsilon}{|\nabla_{\theta_k} \mathcal{L}(\mathcal{D})| + \epsilon}, \quad (8)$$

where ϵ is a small Laplace smoothing constant. Here the second equality follows from the relationship between $\nabla_{\theta_k} \mathcal{L}(\mathcal{D}_{forget})$, $\nabla_{\theta_k} \mathcal{L}(\mathcal{D})$, and $\nabla_{\theta_k} \mathcal{L}(\mathcal{D}_{new})$ that we derived earlier in this section. A small knowledge value characterizes a parameter that is knowledgeable about \mathcal{D}_{forget} , so by selectively re-initializing all parameters θ_k if $\text{QUANTILE}_{KV}(KV_{\theta_k}) \leq \alpha$ (α is a hyperparameter), we can remove the influence of the parameters uniquely responsible for encoding information about these data. This thinking extends on lines of work in gradient-based input saliency maps (Smilkov et al., 2017) and saliency unlearning by Fan et al. (2023). We explore and compare other methods of identifying knowledgeable weights in Appendix A.5.2. This informs Step (4) in Figure 1.

3.3 THE RELOAD ALGORITHM

Based on this intuition, our RELOAD algorithm contains the following steps. (1) Cache the gradients $\nabla_{\theta} \mathcal{L}(\mathcal{D})$ at the end of training. (2) Compute $\nabla_{\theta} \mathcal{L}(\mathcal{D}_{new})$. (3) Perform one step of gradient ascent in the direction of $\nabla_{\theta} (\mathcal{L}(\mathcal{D}) - \mathcal{L}(\mathcal{D}_{new}))$. (4) Re-initialize all parameters θ_k that are smaller than the α -QUANTILE of knowledge values. Finally, (5) fine-tune until convergence on $\mathcal{L}(\mathcal{D}_{new})$. A formal description is shown in Algorithm 1. A software implementation can be found [here](#).

Algorithm 1 The RELOAD Algorithm for Blind Unlearning and Remedial Learning

```

1: Input:  $\mathcal{M}^{(\theta^*)}$ , cached  $\nabla_{\theta} \mathcal{L}(\mathcal{D})$ ,  $\mathcal{D}_{new}$ 
2: Parameters:  $\eta_p$ : priming step learning rate,  $\epsilon$ : noise parameter,  $\alpha$ : reset proportion
3: Output: Trained model approximating  $\mathcal{M}^{(\theta^{\sim})}$ 
4:
5: procedure RELOAD( $\mathcal{M}^{(\theta^*)}$ ,  $\nabla_{\theta} \mathcal{L}(\mathcal{D})$ ,  $\mathcal{M}^{(\theta^*)}$ ,  $\mathcal{D}_{new}$ )
6:    $\theta' \leftarrow \theta^* + \eta_p \nabla_{\theta} (\mathcal{L}(\mathcal{D}) - \mathcal{L}(\mathcal{D}_{new}))$  ▷ Step (2–3) (Fig. 1)
7:    $KV \leftarrow \left\{ \frac{|\nabla_{\theta_k} \mathcal{L}(\mathcal{D}) - \nabla_{\theta_k} \mathcal{L}(\mathcal{D}_{new})| + \epsilon}{|\nabla_{\theta_k} \mathcal{L}(\mathcal{D})| + \epsilon} \right\}_{\theta_k \in \theta}$  ▷ Step (3) (Fig. 1)
8:   for  $\theta_k \in \theta'$  do
9:     if  $\text{QUANTILE}_{KV}(KV_{\theta_k}) \leq \alpha$  then
10:       $\theta'_k \leftarrow \text{INITIALIZE}(\cdot)$  ▷ Step (4) (Fig. 1)
11:    end if
12:  end for
13:  Train  $\mathcal{M}^{(\theta')}$  to convergence on  $\mathcal{D}_{new}$  ▷ Step (5) (Fig. 1)
14: end procedure

```

4 RESULTS AND ANALYSIS

4.1 METHODOLOGICAL INTROSPECTION

Figure 2 introspects on the selected feature maps of a ResNet-18 model when using RELOAD to unlearn the class “8” from the SVHN dataset. The experiment demonstrates the importance of the re-initialization step (Step (4) in Figure 1), as even after a single ascent step, the model still finds “8” to be the most probable class. It is only after the salient weights are identified and re-initialized that the model emits a lower-entropy distribution classifying the digit as a “2”. This suggests that the primary utility of the ascent step in our algorithm is in amending the representations of \mathcal{D}_{forget} in the later layers of the network, while the salient weight re-initialization updates also modify the representations produced by earlier layers. The findings of this experiment present empirical confirmation of the intuition used to develop the algorithm (Section 3.2).

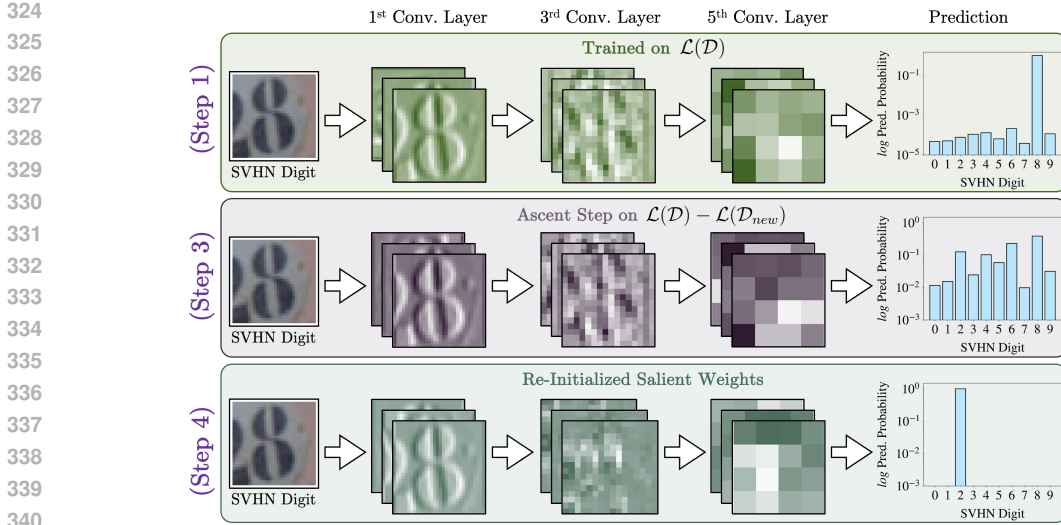


Figure 2: Introspecting on selected feature maps of a ResNet-18 model when using RELOAD to unlearn the class “8”. For brevity, we selected the first channel from each feature map for the sake of visualization, though the patterns we identify appear to hold more broadly across channels. (Top) The feature maps (activations) of the first, third, and fifth convolutional layers after the model was initially trained on \mathcal{D} (Step (1) in Figure 1). (Middle) These same feature masks after the ascent step has been applied (Step (3) in Figure 1). Observe how the activations of the model remain largely unchanged, although the logits represent a considerably more uniform distribution over the digits. (Bottom) These same features masks after the salient weights have been identified and re-initialized (Step (4) in Figure 1). Observe that the activation of the first convolutional layer is largely unchanged – this is expected, as the earlier layers of the network correspond to broad feature detectors (Zeiler & Fergus, 2014) that may be less unique to the features of any particular class in this data. Notice, however, that the feature map of the third convolutional layer is substantially different from that of the previous two stages (it no longer features a hazy “8”), and that the network now emits a significantly lower-entropy distribution predicting the input image as a “2”.

4.2 UNLEARNING EXPERIMENTS

Baselines. We compare RELOAD against baseline approaches of gradient ascent (GA) (Thudi et al., 2022), fine-tuning on \mathcal{D}_{new} (FT) (Warnecke et al., 2023), Selective Synaptic Dampening (SSD) (Foster et al., 2023), SCalable Remembering and Unlearning unBound (SCRUB) (Kurmanji et al., 2023), Catastrophically Forgetting the last k layers (CF- k) (Goel et al., 2022), Exact-Unlearning the last k layers (EU- k) (Goel et al., 2022), SalUn (Fan et al., 2023), and Fisher forgetting (Fisher) (Golatkar et al., 2020). Of these baselines, the requirements of FT, CF- k , EU- k , and Fisher satisfy the blind unlearning setting, whereas the others require direct access to \mathcal{D}_{forget} .

Evaluation. As the goal of approximate unlearning is to produce a learned model that mimics the behaviour of $\mathcal{M}^{(\theta^{\sim})}$, we employ several evaluation statistics to measure the similarity in performance between our learned model and a version of $\mathcal{M}^{(\theta^{\sim})}$ that we naively train from scratch. The accuracy on \mathcal{D}_{new} (NA, \uparrow) measures how well each learned model fits to the new data. The difference in accuracy on \mathcal{D}_{forget} (Δ FA, \downarrow) measures the difference in accuracy between our learned model and $\mathcal{M}^{(\theta^{\sim})}$ on \mathcal{D}_{forget} , while the difference in error on \mathcal{D}_{forget} (Δ FE, \downarrow) measures the difference in (cross-entropy) loss between our learned model and $\mathcal{M}^{(\theta^{\sim})}$ on the \mathcal{D}_{forget} . The difference in success rates of a membership inference attack on \mathcal{D}_{forget} (Δ FMIA, \downarrow) measures the ability of the inference attack from Shokri et al. (2017) to identify members of \mathcal{D}_{forget} in the training data of each learned model, compared to the baseline rate of identification on $\mathcal{M}^{(\theta^{\sim})}$. We also report the AUC of the MIA attack model (Δ AUC, \downarrow). The symmetric KL-divergence on \mathcal{D}_{new} (NSKL, \downarrow) measures the dissimilarity in the logits produced by our learned model and $\mathcal{M}^{(\theta^{\sim})}$ on \mathcal{D}_{new} , while the symmetric KL-divergence on \mathcal{D}_{forget} (FSKL, \downarrow) measures the dissimilarity in the logits produced by our learned model and $\mathcal{M}^{(\theta^{\sim})}$ on \mathcal{D}_{forget} . Cost (\downarrow) measures the computational cost of the method, and is the ratio of time to run the method to the time to naively train $\mathcal{M}^{(\theta^{\sim})}$.

RELOAD unlearns randomly-selected samples. In this experiment, we randomly assign 10% of the training data samples to \mathcal{D}_{forget} , to showcase how well each method can unlearn arbitrary training samples. The results of this experiment are shown in Table 1. Observe that RELOAD achieves the highest NA, while maintaining the lowest Δ FA, Δ FE, Δ FMIA, and FSKL of all approaches. This suggests that RELOAD successfully approximates $\mathcal{M}^{(\theta^{\sim})}$ better than the baselines. That fine-tuning achieves a lower NSKL than RELOAD is hardly surprising, as NSKL measures dissimilarity in logits on \mathcal{D}_{new} , and fine-tuning adjusts a converged model $\mathcal{M}^{(\theta^*)}$ to fit a subset of its original task. Similarly, the computational cost of RELOAD, though similar to many baselines, is considerably greater than either SSD or gradient ascent. The results in Table 1 are produced using a ResNet-18 model on CIFAR-100; additional results with different models and datasets, and on randomly assigning 30% of training data samples to \mathcal{D}_{forget} can be found in Appendix A.6.

Method	NA (\uparrow)	Δ FA (\downarrow)	Δ FE (\downarrow)	Δ FMIA (\downarrow)	Cost (\downarrow)	NSKL (\downarrow)	FSKL (\downarrow)
Retrain	99.98 \pm 0.01	74.89 \pm 2.03	1.06 \pm 0.13	0.63 \pm 0.20	1.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00
GA	93.81 \pm 0.75	18.77 \pm 2.43	0.95 \pm 0.14	0.21 \pm 0.06	0.00\pm0.00	0.29 \pm 0.09	2.62 \pm 0.05
FT	96.00 \pm 0.12	16.46 \pm 2.47	0.89 \pm 0.14	0.19 \pm 0.08	0.27 \pm 0.00	0.03\pm0.01	2.11 \pm 0.06
SSD	1.01 \pm 0.02	74.17 \pm 2.04	4.19 \pm 0.59	0.15 \pm 0.21	0.01 \pm 0.00	14.90 \pm 1.24	11.81 \pm 1.24
SCRUB	93.76 \pm 0.74	18.85 \pm 2.39	0.95 \pm 0.14	0.20 \pm 0.06	0.02 \pm 0.00	0.29 \pm 0.09	2.63 \pm 0.06
CF- <i>k</i>	94.75 \pm 0.41	18.01 \pm 2.60	0.94 \pm 0.14	0.20 \pm 0.06	0.21 \pm 0.00	0.14 \pm 0.03	2.47 \pm 0.07
EU- <i>k</i>	94.32 \pm 0.49	17.93 \pm 2.55	0.94 \pm 0.14	0.20 \pm 0.06	0.21 \pm 0.00	0.19 \pm 0.05	2.33 \pm 0.05
SalUn	99.06 \pm 0.22	13.14 \pm 2.53	0.11 \pm 0.09	7.39 \pm 2.60	0.16 \pm 0.00	0.06 \pm 0.02	0.55\pm0.04
Fisher	97.76 \pm 0.78	22.99 \pm 2.30	0.95 \pm 0.14	7.27 \pm 2.48	1.78 \pm 0.04	0.07 \pm 0.02	0.56 \pm 0.04
RELOAD	99.56\pm0.11	0.30\pm0.50	0.04\pm0.02	0.01\pm0.01	0.12 \pm 0.01	0.15 \pm 0.03	1.23 \pm 0.11

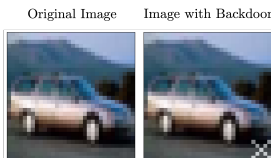
Table 1: **10% Random Forgetting on CIFAR-100 (ResNet-18).** The top row presents the value of $\mathcal{M}^{(\theta^{\sim})}$ on each metric. Subsequent rows for Δ FA (\downarrow), Δ FE (\downarrow), and Δ FMIA (\downarrow) present the absolute difference in the value of the corresponding method on this metric to the value of $\mathcal{M}^{(\theta^{\sim})}$ on the metric. These results show that RELOAD outperforms all the baselines on NA, Δ FA, Δ FE, Δ FMIA, and FSKL by large margins. RELOAD performs competitively on the NSKL metric, outperformed by FT and CF-*k*. RELOAD incurs a higher computational cost than most baselines, but is cheaper than FT, CF-*k*, and EU-*k*.

Method	NA (\uparrow)	FA (Δ^{\downarrow})	FE (Δ^{\downarrow})	FMIA (Δ^{\downarrow})	Cost (\downarrow)	NSKL (\downarrow)	FSKL (\downarrow)
Retrain	99.99 \pm 0.00	95.12 \pm 0.23	0.20 \pm 0.01	0.50 \pm 0.00	1.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00
GA	99.57 \pm 0.02	4.37 \pm 0.25	0.17 \pm 0.01	0.05 \pm 0.01	0.00\pm0.00	0.05 \pm 0.00	0.52 \pm 0.02
FT	99.99\pm0.00	4.33 \pm 0.22	0.17 \pm 0.01	0.04 \pm 0.01	0.27 \pm 0.00	0.00\pm0.00	0.43 \pm 0.02
SSD	12.75 \pm 4.69	82.52 \pm 4.73	2.12 \pm 0.06	0.01 \pm 0.01	0.01 \pm 0.00	8.55 \pm 0.13	7.88 \pm 0.12
SCRUB	99.79 \pm 0.01	4.44 \pm 0.26	0.18 \pm 0.01	0.05 \pm 0.01	0.03 \pm 0.00	0.03 \pm 0.00	0.50 \pm 0.02
CF- <i>k</i>	99.76 \pm 0.01	4.47 \pm 0.24	0.18 \pm 0.01	0.05 \pm 0.01	0.23 \pm 0.02	0.03 \pm 0.00	0.50 \pm 0.02
EU- <i>k</i>	99.63 \pm 0.01	4.46 \pm 0.25	0.18 \pm 0.01	0.05 \pm 0.01	0.23 \pm 0.02	0.05 \pm 0.00	0.47 \pm 0.02
SalUn	99.90 \pm 0.04	3.14 \pm 1.00	0.13 \pm 0.03	0.04 \pm 0.02	0.17 \pm 0.00	0.03 \pm 0.00	0.50 \pm 0.02
Fisher	99.57 \pm 0.02	0.09\pm0.05	0.00\pm0.00	0.01 \pm 0.00	2.17 \pm 0.04	0.05 \pm 0.00	0.47 \pm 0.02
RELOAD	99.68 \pm 0.17	0.25 \pm 0.21	0.01 \pm 0.01	0.00\pm0.00	0.12 \pm 0.01	0.06 \pm 0.02	0.21\pm0.02

Table 2: **100 In Class Random Forgetting on SVHN (ResNet-18)**

\uparrow : the goal is to have as high of a value as possible, Δ^{\downarrow} : the value in the table is the difference between the result of the unlearning method and retraining (top row) on the metric and the goal is to have a low difference, \downarrow : the goal is to have as low of a value as possible. The top row presents the value of $\mathcal{M}^{(\theta^{\sim})}$ on each metric. Subsequent rows for Δ FA (\downarrow), Δ FE (\downarrow), and Δ FMIA (\downarrow) present the absolute difference in the value of the corresponding method on this metric to the value of $\mathcal{M}^{(\theta^{\sim})}$ on the metric. These results show that RELOAD outperforms all the baselines on Δ FA, Δ FE, Δ FMIA, and NSKL by large margins. RELOAD performs competitively on NA and FSKL but is outperformed by FT. RELOAD also incurs a higher computational cost than the other baselines.

RELOAD efficiently unlearns correlated samples. We next randomly assign 100 samples from a single class of the training data to \mathcal{D}_{forget} , to evaluate how well each method can unlearn arbitrary but related training samples. The results of this experiment are shown in Table 2.



RELOAD achieves the lowest Δ FMIA and FSKL of all approaches and very close to the lowest Δ FA, Δ FE, and NSKL of all approaches, suggesting that again, RELOAD learns to closely approximate $\mathcal{M}^{(\theta^{\sim})}$ in this setting. RELOAD is marginally outperformed

Figure 3: The “Cross Pattern Backdoor” inserts the above pattern (right) in all images from

432 by Fisher in these settings, but is far more realistic, as Fisher over
 433 twice as much time as retraining. Although RELOAD achieves an
 434 NA competitive with that of most baselines, naive gradient ascent,
 435 CF- k , and EU- k yield a marginally higher NA; this is surprising for
 436 gradient ascent as it typically yields lower NA values. This can be
 437 attributed to the small number of unlearning samples; optimizing
 438 to maximize the loss on these samples does not provide much of
 439 a gradient update. CF- k and EU- k both make few parameter up-
 440 dates to $\mathcal{M}^{(\theta^*)}$, which leads to a high NA but poor performance on
 441 unlearning metrics like Δ FA and Δ FE.

442 4.3 REMEDIAL LEARNING EXPERIMENTS

443 **Baselines.** The remedial unlearning setting admits different baselines than the unlearning setting.
 444 “Original” represents a baseline model trained on \mathcal{D} , and “Retrain” represents a baseline model
 445 trained directly on \mathcal{D}_{new} . Then, because gradient ascent does not directly translate to the task of
 446 remedial learning (because there is no “forget set” on which to ascent) we introduce two variants
 447 of gradient ascent to serve as baselines. Gradient Ascent Relearn (GAR) performs 10 epochs of
 448 gradient ascent on \mathcal{D} , followed by 10 epochs of gradient descent on \mathcal{D}_{new} . Gradient Difference
 449 Ascent (GRDA) calculates $\nabla\mathcal{L}(\mathcal{D}) - \nabla\mathcal{L}(\mathcal{D}_{new})$ on each step and performs a gradient update in the
 450 opposite direction, fitting \mathcal{D}_{new} . It performs 10 epochs of such updates.

451 **Metrics.** We evaluate remedial learning as follows. The accuracy on a held-out test split of \mathcal{D}_{new}
 452 (Acc. $\mathcal{D}_{new}^{(test)}$, \uparrow) represents how well the model fits the distribution of \mathcal{D}_{new} (without the backdoor;
 453 see next section). The accuracy on a held-out test split of \mathcal{D}_{new} with backdoors added to each
 454 instance (Acc. $\mathcal{D}_{new}^{(test,\S)}$) measures the reliance of the model on the backdoor pattern in classification.
 455 Specifically, a low accuracy on $\mathcal{D}_{new}^{(test,\S)}$ suggests that the model is over-reliant on the backdoor
 456 pattern that was injected into its training data.

457 **RELOAD corrects erroneous data (re-**

458 **moving shortcuts).** In this setting, we select 2 classes from the data (here, CIFAR-10) and inject cross-patterns into the corners of their training samples to construct \mathcal{D} . An example of this backdoor can be seen in Figure 3. The inclusion of this backdoor influences a model trained on this dataset to rely on the cross-patterns as strong indicators of class membership. To construct \mathcal{D}_{new} we then replace the cross-patterned samples with their original instances, removing the backdoor. The goal of remedial learning here is to un-learn the model’s reliance on the backdoor, and re-learn the salient representations needed to accurately predict on the affected classes. The results of this experiment are shown in Appendix A.9. Observe that the effect of this backdoor attack produces a trained model (Original; trained *with* the backdoor) with poor performance on $\mathcal{D}_{new}^{(test,\S)}$, because the model learned to treat the backdoor pattern as a strong indicator of class membership for certain classes. Further, notice that RELOAD successfully remedies this vulnerability, achieving the highest accuracy (aside from Retrain; retrained from scratch *without* the backdoor) on $\mathcal{D}_{new}^{(test)}$, and the highest accuracy of all models on $\mathcal{D}_{new}^{(test,\S)}$. This suggests that RELOAD is capable of efficiently correcting the predictive behaviour of a model trained on erroneous data.

Method	Acc. $\mathcal{D}_{new}^{(test)}$ (\uparrow)	Acc. $\mathcal{D}_{new}^{(test,\S)}$ (\uparrow)	Cost (\downarrow)
Original	82.68 \pm 0.45	19.81 \pm 0.03	N/A
Retrain	92.48 \pm 0.00	91.90 \pm 0.00	1.00 \pm 0.00
GAR	57.29 \pm 34.88	56.54 \pm 34.12	0.08 \pm 0.01
GRDA	62.87 \pm 28.47	62.34 \pm 27.80	0.05 \pm 0.00
FT	86.87 \pm 4.39	86.50 \pm 4.14	0.37 \pm 0.02
SSD	30.25 \pm 22.90	23.94 \pm 13.43	0.01\pm0.00
SCRUB	12.43 \pm 3.45	12.42 \pm 3.44	0.04 \pm 0.01
CF- k	66.56 \pm 24.27	66.29 \pm 23.80	0.29 \pm 0.03
EU- k	66.75 \pm 24.08	66.41 \pm 23.63	0.29 \pm 0.03
RELOAD	90.81\pm0.99	90.51\pm0.82	0.08 \pm 0.06

459 Table 3: **Cross Pattern Backdoor Removal on CIFAR-10 (ResNet-18).** \uparrow : the goal is to have as high of a value as possible, \downarrow : the goal is to have as low of a value as possible. These results show that RELOAD outperforms all baselines on Acc. $\mathcal{D}_{new}^{(test)}$ and Acc. $\mathcal{D}_{new}^{(test,\S)}$. The small differences between these accuracy values for RELOAD indicate that it successfully removed the influence of the backdoor pattern. RELOAD incurs a higher computational cost than most baselines, but is cheaper than FT, CF- k , and EU- k .

486 5 DISCUSSION
487

488 This work introduces the setting of *blind unlearning*, machine unlearning without direct access to the
489 “forget set”. This setting allows for improved privacy procedures in practical settings, by enabling
490 the immediate deletion of data when an unlearning request is received rather than retaining the data
491 for the purpose of downstream unlearning. Our method, RELOAD, combines insights from gradient-
492 based unlearning (to remove top-level information from all parameters) with selective parameter
493 re-initialization. The blind setting ensures that as long as practitioners store the last step gradients
494 of their model on the training set, they have the capacity to unlearn data when it is removed from
495 their system. We recommend that future work study the performance of RELOAD at larger scales,
496 such as those presented by modern large language models (Achiam et al., 2023), and investigate the
497 utility of other choices for $\mathcal{I}_{\mathcal{D}}$ beyond the cached gradients used in RELOAD.

498 Despite operating in the blind setting, RELOAD outperforms benchmark machine unlearning algo-
499 rithms that enjoy direct access to \mathcal{D}_{forget} , suggesting that it is an empirically effective unlearning
500 algorithm. However, RELOAD admits a modest tradeoff between computational efficiency and per-
501 formance in this regime. We finally observe that machine unlearning represents a special case of
502 *remedial learning*, a setting that is especially important for efficiently correcting errors in the train-
503 ing data used to train models. RELOAD remains an efficient, performant method in this regime,
504 suggesting that our work may contain generalizable insights about about learning to fit arbitrary
505 downstream transformations of data.
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

REFERENCES

- 540
541
542 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-
543 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
544 report. *arXiv preprint arXiv:2303.08774*, 2023.
- 545 Lucas Bourtole, Varun Chandrasekaran, Christopher A. Choquette-Choo, Hengrui Jia, Adelin
546 Travers, Baiwu Zhang, David Lie, and Nicolas Papernot. Machine unlearning. 12 2019. URL
547 <http://arxiv.org/abs/1912.03817>.
- 548 Yinzhi Cao and Junfeng Yang. Towards making systems forget with machine unlearning. In *2015*
549 *IEEE Symposium on Security and Privacy*, pp. 463–480, 2015. doi: 10.1109/SP.2015.35.
- 550
551 Asic Chen, Ruian Ian Shi, Xiang Gao, Ricardo Baptista, and Rahul G Krishnan. Structured neural
552 networks for density estimation and causal inference. *Advances in Neural Information Processing*
553 *Systems*, 36, 2024.
- 554 Min Chen, Weizhuo Gao, Gaoyang Liu, Kai Peng, and Chen Wang. Boundary unlearning: Rapid
555 forgetting of deep networks via shifting the decision boundary. *2023 IEEE/CVF Conference*
556 *on Computer Vision and Pattern Recognition (CVPR)*, pp. 7766–7775, 2023. URL <https://api.semanticscholar.org/CorpusID:257636742>.
- 557
558 Eli Chien, Chao Pan, and Olgica Milenkovic. Certified graph unlearning. *ArXiv*, abs/2206.09140,
559 2022. URL <https://api.semanticscholar.org/CorpusID:249890116>.
- 560
561 Vikram S Chundawat, Ayush K Tarun, Murari Mandal, and Mohan S. Kankanhalli. Can bad
562 teaching induce forgetting? unlearning in deep networks using an incompetent teacher. *ArXiv*,
563 abs/2205.08096, 2022. URL <https://api.semanticscholar.org/CorpusID:248834527>.
- 564
565 Marco Cotogni, Jacopo Bonato, Luigi Sabetta, Francesco Pelosin, and Alessandro Nicolosi. Duck:
566 Distance-based unlearning via centroid kinematics. *ArXiv*, abs/2312.02052, 2023. URL <https://api.semanticscholar.org/CorpusID:265609937>.
- 567
568 European Parliament and Council of the European Union. Regulation (EU) 2016/679 of the Euro-
569 pean Parliament and of the Council. URL [https://data.europa.eu/eli/reg/2016/](https://data.europa.eu/eli/reg/2016/679/oj)
570 [679/oj](https://data.europa.eu/eli/reg/2016/679/oj).
- 571
572 Chongyu Fan, Jiancheng Liu, Yihua Zhang, Eric Wong, Dennis Wei, and Sijia Liu. Salun: Em-
573 powering machine unlearning via gradient-based weight saliency in both image classification and
574 generation, 2023. URL <https://arxiv.org/abs/2310.12508>.
- 575
576 Farzaneh S. Fard, Paul Hollensen, Stuart Mcilory, and Thomas Trappenberg. Impact of biased
577 mislabeling on learning with deep networks. In *2017 International Joint Conference on Neural*
Networks (IJCNN), pp. 2652–2657, 2017. doi: 10.1109/IJCNN.2017.7966180.
- 578
579 Jack Foster, Stefan Schoepf, and Alexandra Brintrup. Fast machine unlearning without retraining
580 through selective synaptic dampening. *ArXiv*, abs/2308.07707, 2023. URL <https://api.semanticscholar.org/CorpusID:260900355>.
- 581
582 Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural
583 networks. *arXiv preprint arXiv:1803.03635*, 2018.
- 584
585 Kai-Xin Gao, Xiaolei Liu, Zheng-Hai Huang, Min Wang, Shuangling Wang, Zidong Wang, Dachuan
586 Xu, and F. Yu. Eigenvalue-corrected natural gradient based on a new approximation. *Asia Pac.*
587 *J. Oper. Res.*, 40:2340005:1–2340005:18, 2020. URL <https://api.semanticscholar.org/CorpusID:227210187>.
- 588
589 Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel,
590 Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature*
591 *Machine Intelligence*, 2(11):665–673, 2020.
- 592
593 Shashwat Goel, Ameya Prabhu, and Ponnurangam Kumaraguru. Evaluating inexact unlearn-
ing requires revisiting forgetting. *ArXiv*, abs/2201.06640, 2022. URL <https://api.semanticscholar.org/CorpusID:246015741>.

- 594 Aditya Golatkar, Alessandro Achille, and Stefano Soatto. Eternal sunshine of the spotless net: Se-
595 lective forgetting in deep networks. In *2020 IEEE/CVF Conference on Computer Vision and*
596 *Pattern Recognition (CVPR)*. IEEE, June 2020. doi: 10.1109/cvpr42600.2020.00932. URL
597 <http://dx.doi.org/10.1109/CVPR42600.2020.00932>.
- 598
599 Laura Graves, Vineel Nagisetty, and Vijay Ganesh. Amnesiac machine learning. In *Proceedings of*
600 *the AAAI Conference on Artificial Intelligence*, volume 35, pp. 11516–11524, 2021.
- 601 Chuan Guo, Tom Goldstein, Awni Y. Hannun, and Laurens van der Maaten. Certified data re-
602 moval from machine learning models. *ArXiv*, abs/1911.03030, 2019. URL [https://api.](https://api.semanticscholar.org/CorpusID:207847600)
603 [semanticscholar.org/CorpusID:207847600](https://api.semanticscholar.org/CorpusID:207847600).
- 604 B. Hassibi, D.G. Stork, and G.J. Wolff. Optimal brain surgeon and general network pruning. In
605 *IEEE International Conference on Neural Networks*, pp. 293–299 vol.1, 1993. doi: 10.1109/
606 ICNN.1993.298572.
- 607
608 Tien Ho-Phuoc. Cifar10 to compare visual recognition performance between deep neural networks
609 and humans. *ArXiv*, abs/1811.07270, 2018. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:53721934)
610 [CorpusID:53721934](https://api.semanticscholar.org/CorpusID:53721934).
- 611 Alex Krizhevsky. Learning multiple layers of features from tiny images. *University of Toronto*, 05
612 2012.
- 613
614 Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-100 (canadian institute for advanced
615 research). URL <http://www.cs.toronto.edu/~kriz/cifar.html>.
- 616 Meghdad Kurmanji, Peter Triantafillou, Jamie Hayes, and Eleni Triantafillou. Towards unbounded
617 machine unlearning. In *Thirty-seventh Conference on Neural Information Processing Systems*,
618 2023. URL <https://openreview.net/forum?id=OveBaTtUAT>.
- 619
620 Jing Lin, Ryan S. Luley, and Kaiqi Xiong. Active learning under malicious mislabeling and poi-
621 soning attacks. *2021 IEEE Global Communications Conference (GLOBECOM)*, pp. 1–6, 2021.
622 URL <https://api.semanticscholar.org/CorpusID:230437562>.
- 623 James Martens and Roger Baker Grosse. Optimizing neural networks with kronecker-factored ap-
624 proximate curvature. In *International Conference on Machine Learning*, 2015. URL [https:](https://api.semanticscholar.org/CorpusID:11480464)
625 [//api.semanticscholar.org/CorpusID:11480464](https://api.semanticscholar.org/CorpusID:11480464).
- 626 Yuval Netzer, Tao Wang, Adam Coates, A. Bissacco, Bo Wu, and A. Ng. Reading digits in natural
627 images with unsupervised feature learning. 2011. URL [https://api.semanticscholar.](https://api.semanticscholar.org/CorpusID:16852518)
628 [org/CorpusID:16852518](https://api.semanticscholar.org/CorpusID:16852518).
- 629
630 Jonas Pfeiffer, Sebastian Ruder, Ivan Vulić, and Edoardo Ponti. Modular deep learning. *Transactions*
631 *on Machine Learning Research*, 2023. ISSN 2835-8856. URL [https://openreview.net/](https://openreview.net/forum?id=z9EkXfvxta)
632 [forum?id=z9EkXfvxta](https://openreview.net/forum?id=z9EkXfvxta). Survey Certification.
- 633 Nathaniel Rodriguez, E. Izquierdo, and Yong-Yeol Ahn. Optimal modularity and memory capacity
634 of neural reservoirs 1, 2019.
- 635
636 Thanveer Shaik, Xiaohui Tao, Haoran Xie, Lin Li, Xiaofeng Zhu, and Qing Li. Exploring the
637 landscape of machine unlearning: A comprehensive survey and taxonomy, 2024.
- 638 R. Shokri, M. Stronati, C. Song, and V. Shmatikov. Membership inference attacks against machine
639 learning models. In *2017 IEEE Symposium on Security and Privacy (SP)*, pp. 3–18, Los Alamitos,
640 CA, USA, may 2017. IEEE Computer Society. doi: 10.1109/SP.2017.41. URL [https://doi.](https://doi.ieeeecomputersociety.org/10.1109/SP.2017.41)
641 [ieeecomputersociety.org/10.1109/SP.2017.41](https://doi.ieeeecomputersociety.org/10.1109/SP.2017.41).
- 642 Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda B. Viégas, and Martin Wattenberg. Smooth-
643 grad: removing noise by adding noise. *ArXiv*, abs/1706.03825, 2017. URL [https://api.](https://api.semanticscholar.org/CorpusID:11695878)
644 [semanticscholar.org/CorpusID:11695878](https://api.semanticscholar.org/CorpusID:11695878).
- 645
646 Ayush K. Tarun, Vikram S. Chundawat, Murari Mandal, and Mohan Kankanhalli. Fast yet effective
647 machine unlearning. *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–10,
2023. doi: 10.1109/TNNLS.2023.3266233.

648 A. Thudi, G. Deza, V. Chandrasekaran, and N. Papernot. Unrolling sgd: Understanding factors
649 influencing machine unlearning. In *2022 IEEE 7th European Symposium on Security and Privacy*
650 *(EuroS&P)*, pp. 303–319, Los Alamitos, CA, USA, jun 2022. IEEE Computer Society. doi:
651 10.1109/EuroSP53844.2022.00027. URL [https://doi.ieeecomputersociety.org/
652 10.1109/EuroSP53844.2022.00027](https://doi.ieeecomputersociety.org/10.1109/EuroSP53844.2022.00027).

653 Alexander Warnecke, Lukas Pirch, Christian Wressnegger), and Konrad Rieck. Machine unlearning
654 of features and labels. In *Proceedings 2023 Network and Distributed System Security Symposium,*
655 *NDSS 2023*. Internet Society, 2023. doi: 10.14722/ndss.2023.23087. URL [http://dx.doi.
656 org/10.14722/ndss.2023.23087](http://dx.doi.org/10.14722/ndss.2023.23087).

657 Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In
658 *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12,*
659 *2014, Proceedings, Part I 13*, pp. 818–833. Springer, 2014.

660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701