
TACKLING FEDERATED UNLEARNING AS A PARAMETER ESTIMATION PROBLEM

Antonio Balordi

CASD - Italian Defense University
Rome, Italy
antonio.balordi@unicas.it

Lorenzo Manini

DIFA - University of Bologna
Bologna, Italy
lorenzo.manini@studio.unibo.it

Fabio Stella

Department of Informatics, Systems and Communication,
University of Milano-Bicocca,
Milan, Italy
fabio.stella@unimib.it

Alessio Merlo

CASD - Italian Defense University
Rome, Italy
alessio.merlo@unicas.it

ABSTRACT

Privacy regulations require the erasure of data from deep learning models. This is a significant challenge that is amplified in Federated Learning, where data remains on clients, making full retraining or coordinated updates often infeasible. This work introduces an efficient Federated Unlearning framework based on information theory, modeling leakage as a parameter estimation problem. Our method uses second-order Hessian information to identify and selectively reset only the parameters most sensitive to the data being forgotten, followed by minimal federated retraining. This model-agnostic approach supports categorical and client unlearning without requiring server access to raw client data after initial information aggregation. Evaluations on benchmark datasets demonstrate strong privacy (MIA success near random, categorical knowledge erased) and high performance (Normalized Accuracy against re-trained benchmarks of ≈ 0.9), while aiming for increased efficiency over complete retraining. Furthermore, in a targeted backdoor attack scenario, our framework effectively neutralizes the malicious trigger, restoring model integrity. This offers a practical solution for data forgetting in FL.

Keywords Federated Unlearning · Parameter estimation · Right To Be Forgotten · Data Privacy · Hessian computation.

1 Introduction

The rapid adoption of deep learning has led to models trained on massive datasets that achieve outstanding results across domains. However, such outcomes come at the cost of privacy: Neural networks (NN) are known to memorize training data, potentially leaking sensitive information ([1], [2]). This privacy concern has led to *machine unlearning* gaining momentum, which generally refers to a set of techniques aimed at allowing organizations to meet the proper privacy constraints by selectively removing the influence of specific data from trained models. This capability has recently become crucial as privacy regulations, such as GDPR ([3], [4]) or CCPA ([5], [6]), provide the user with the “*right to be forgotten*”. As a consequence, clients can request the complete removal of their data in any training scenario.

This work will focus on Federated Learning (FL), a learning technique that allows users to collectively obtain the benefits of shared models without the need to share the data [7]. Despite this privacy-preserving design, the need to retroactively remove a user’s data contribution presents a significant challenge, leading to the development of the field of *Federated Unlearning (FU)*. Here, the *gold standard* refers to the practice of *full retraining*, that is, the model is discarded and a new one is trained from scratch without the removed client; this approach grants an *exact unlearning* ([8], [9]). However, it often leads to significant communication overhead, inefficient resource utilization, and high computational costs, rendering it impractical for many real-world federated learning (FL) deployments ([10], [11]).

Consequently, applying exact unlearning through full retraining reveals that it is often unfeasible in FL scenarios with high client turnover or when dealing with large models and datasets [12]. Therefore, there is a strong need for FU methods that can approximate the guarantees provided by the gold standard while operating within feasible resource constraints and avoiding the need to retrain from scratch.

Contributions. This paper introduces a new federated unlearning framework that conceptualizes information leakage as a *parameter estimation problem*. This information-theoretic perspective enables the development of a high-performing unlearning method that selectively resets and retrains (minimally, for one or a few epochs) only those parameters identified as most informative about the data to be unlearned, utilizing second-order information (Hessian diagonals) computed post-training. Our approach is broadly applicable across diverse network architectures, supports arbitrary dataset unlearning, and integrates seamlessly into the federated setting, enabling server-side unlearning with no client data access and retraining overhead.

Paper Structure. The remainder of the paper is organized as follows. Section 2 reviews relevant previous work on federated unlearning. Section 3 introduces the federated unlearning problem. Our proposed framework is detailed in Section 4, followed by a description of the experimental setup and datasets in Section 6. We present some numerical results in Section 7 and discuss them in Section 8. Section 9 concludes the paper, summarizing key findings and describing future directions.

2 Related Work

Machine Unlearning (MU) addresses the need to remove the influence of specific data from a trained model, a capability required to ensure privacy and comply with regulatory requirements [13]. In centralized settings, extensive research has explored various techniques. Among these, methods such as SISA (sharded retraining) [11], which focuses on efficient exact removal via sharding, and influence functions [14], which offer a theoretical lens on the impact of data, have been extensively explored.

However, translating these MU principles to Federated Learning (FL) contexts, a domain we refer to as Federated Unlearning (FU), introduces new challenges. The direct application of centralized MU methods to FU scenarios is often problematic. For example, techniques such as SISA or influence functions can be complex to implement effectively with deep learning models typically employed in FL, and their unlearning guarantees may not extend to these distributed environments [15]. For SISA, adapting its isolated shard retraining paradigm to the naturally distributed individuals’ data in FL without substantial re-coordination or compromising global model integrity is complex. Influence functions, on the other hand, typically require computationally intensive calculations that are challenging for deep FL models and often need direct data access, which violates FL’s core privacy tenets.

Moreover, FU must inherently work with the complexities of FL itself, including the distributed nature of client data, communication bottlenecks, and increased privacy considerations [16, 17].

Current FU approaches can be broadly categorized as follows, each with significant limitations.

- **Retraining-Based/Approximation:** These methods approximate complete retraining, for example, through knowledge distillation [18], [19]. Although faster, they still require multiple communication rounds between clients participating in the training and significant computation. More critically, in knowledge-distillation setups the *student* often inherits the *teacher’s* membership signals, so a model can ‘pass’ one audit while still retaining information—making the quality of forgetting hard to measure and increasingly fragile as the forget set or model complexity grows [15, 20].
- **Gradient/Update Manipulation:** Several approaches attempt to reverse the impact of unlearned data by manipulating gradients or model updates. For example, [21] proposed methods including gradient ascent on the forgotten client’s data. However, these are often heuristic; forcing misprediction does not guarantee removal of memorized information and can leave the model vulnerable to reconstruction or membership inference attacks since it could still retain encoded information, as shown in [2] and in [22].
- **Parameter Masking/Perturbation:** These techniques alter parameters related to forgotten data, for example, pruning or adding noise [23], [24]. Identifying relevant parameters is non-trivial; currently, most methods typically fail to capture complex dependencies [25]. Determining the extent of the modification requires tuning and risk degradation of the utility [26]. Addressing this, our proposed methodology (Section 4) offers an information-theoretic strategy for the identification and reset of the targeted parameters within this general class of methods.

3 Problem Setup

3.1 Federated Learning and Unlearning

In FL, some clients and a central aggregating server collaboratively train a machine learning model while maintaining the locality of the data, thus obtaining improved privacy and security.

Then, FU aims to achieve one or more of the following objectives [10]:

- **Sample Unlearning:** the removal of a specific subset of data samples.
- **Client Unlearning:** the removal of the entire dataset of a given client. This scenario most directly corresponds to the “*Right to Be Forgotten*” (RTBF) principle.
- **Class Unlearning:** the removal of all data samples belonging to a specific class c , potentially affecting the data contributions of multiple clients.

It should be noted that both Client Unlearning and Class Unlearning can be viewed as specific instances of the Sample Unlearning problem. In each case, the goal is to remove a subset of samples defined by a particular characteristic (i.e., belonging to a specific client or class).

In FU, the unlearning process should provide the same level of privacy as the original federated learning process: the clients’ data should remain private and not be accessible by the server or other clients.

3.2 Challenges in Federated Unlearning

The requirements outlined in Section 3.1 pose a well-defined set of challenges, as effectively outlined in [10]. We revisit key aspects of these challenges to motivate and present the specific features of our proposed Federated Unlearning framework and methodology.

- *Nondeterministic Training Process:* Dynamic client participation in FL complicates the exact reproduction of past training states for unlearning (often leading to discarded knowledge). Although the theoretical framework we present is developed under the assumption of deterministic training, empirical results in Section 7 demonstrate its effectiveness in more realistic scenarios, i.e., nondeterministic scenarios.
- *Inscrutable Model Updates and Finite Memory:* The server’s potential inability to access individual client updates and the impracticality of storing all historical model states restrict methods reliant on reversing past aggregations. To avoid this, our framework leverages the computation of second-order statistics (Hessian diagonals or derived information scores) on the final model state.
- *Inscrutable Local Datasets:* By design in FL, client data remains private and cannot be accessed by the server or other clients for unlearning. This prevents methods that require direct access to the data. Our framework upholds this principle by requiring clients to share only second-order statistics and not the raw data itself. The subsequent retraining phase also adheres to standard FL protocols.
- *Iterative Learning Process:* The iterative nature of FL means that a naive rollback to remove the client’s contribution from a previous round would invalidate subsequent aggregations. Our approach directly modifies the current model state, removing the specific imprint of the unlearned data, followed by a brief reconvergence phase with the remaining knowledge. This setting is detailed in our experiments.
- *Data Heterogeneity:* In an FL setting, data might often be non-IID across clients; however, due to the secrecy of the clients’ datasets, assessing the data heterogeneity is not straightforward. Nevertheless, the data-agnostic model of our framework is well-suited to leverage this, as it inherently adapts to data heterogeneity without requiring prior knowledge of its extent.

4 Parameter Estimation Framework

To characterize the information encoded in each of the network parameters, we consider a scenario in which a hypothetical attacker attempts to infer details about the target dataset (TD) ¹ through interactions with the neural network. We model this attack as a *parameter estimation problem* where the parameter to estimate is whether the target dataset was used during training or not, and the observables are the NN parameters. We then leverage results from

¹A formal definition of the Target Dataset will be provided in Section 5

information theory to define a *Target Information Score (TIS)* that quantifies how much each parameter contributes to the estimability of information about the TD, and in Section 4 we use it to define a procedure to forget the TD.

The attack scenario is defined as follows:

- The attacker has full knowledge about the training procedure. The training procedure is modelled as a function² $\mathcal{T} : \mathcal{D} \rightarrow \mathfrak{P}$ from the space of all possible datasets \mathcal{D} to the space of all possible network parameters \mathfrak{P} .
- The attacker knows the trained network parameters $\tilde{\theta} \in \mathfrak{P}$.
- The attacker is aware that only two scenarios are possible: either the target dataset (TD) was included in the training process, or it was entirely excluded. Formally, let us denote the TD as $\tilde{D}_T \in \mathcal{D}$, then the two scenarios correspond to two possible training datasets
 - $\tilde{D}_1 \in \mathcal{D}$ such that $\tilde{D}_T \subset \tilde{D}_1$.
 - $\tilde{D}_0 \in \mathcal{D}$ such that $\tilde{D}_0 = \tilde{D}_1 \setminus \tilde{D}_T$.

The information the attacker wants to infer is which of the two datasets was used; formally, she thus wants to estimate a binary parameter

$$t = \begin{cases} 0 & \text{if the training dataset was } \tilde{D}_0 \\ 1 & \text{if the training dataset was } \tilde{D}_1 \end{cases} \quad (1)$$

so that we also denote the (parametric) training dataset with \tilde{D}_t .

- The attacker has an "arbitrarily good" subjective probability distribution about the possible training datasets. Precisely, we consider the attacker not to know the exact datasets \tilde{D}_t but that her knowledge provides her with two subjective probability distributions defined on \mathcal{D} and peaked around \tilde{D}_t , as follows

$$\rho_t^{data}(D) = \mathcal{N}(D; \tilde{D}_t, \epsilon \mathbb{1}) \quad t \in \{0, 1\} \quad (2)$$

where \mathcal{N} is a multivariate Gaussian distribution defined on the sample space of each dataset and ϵ quantifies the uncertainty about the datasets³. We thus have that in the limit of vanishing ϵ the prior tends to a Dirac delta centered on the true datasets, i.e., exact knowledge of them. Moreover, from their definitions, it follows that

$$\rho_1^{data}(D) = \mathcal{N}(D; \tilde{D}_0, \epsilon \mathbb{1}) \mathcal{N}(D; \tilde{D}_T, \epsilon \mathbb{1}) \quad (3)$$

$$= \rho_0^{data}(D) \rho_T^{data}(D) \quad (4)$$

since \tilde{D}_1 is the disjoint union of \tilde{D}_0 and \tilde{D}_T , and where we defined the distribution of the TD as

$$\rho_T^{data}(D) = \mathcal{N}(D; \tilde{D}_T, \epsilon \mathbb{1}) \quad (5)$$

From the attacker's point of view, then her lack of knowledge about the precise content of the two possible datasets, quantified by ϵ , leads to a lack of knowledge about the expected trained network parameters in the two possible cases. Specifically, this is described by the fact that since the training dataset D is a random outcome sampled by one of the two probability distributions $\rho_t^{data}(D)$. The trained network parameters are a random variable given by $\theta = \mathcal{T}(D)$. Thus, we will have two probability distributions for the trained parameters defined in \mathfrak{P} :

$$\rho_t^{par}(\theta) \quad t \in \{0, 1\} \quad (6)$$

Then, if both $\rho_0^{par}(\tilde{\theta}) > 0$ and $\rho_1^{par}(\tilde{\theta}) > 0$, the attacker will not be able to infer with complete certainty the value of t . Still, she will only be able to assign a likelihood to the two possibilities. Formally, the procedure the attacker uses to obtain his best guess of the parameter t is an estimator $\hat{t} : \mathfrak{P} \rightarrow \{0, 1\}$.

We now want to study the information-theoretical limits to the efficacy of the attack procedure, that is, to the efficacy of any estimator \hat{t} . Our analysis will ultimately take advantage of theorems that provide a lower bound on the variance of \hat{t} . However, these theorems require the parameter to be estimated to be continuous, which implies having a one-dimensional family of probability distributions. To satisfy this hypothesis, we need a continuous embedding of t , the choice of which is not, in general, unique. To fix ideas, one might picture the embedding as choosing any smooth interpolation from \tilde{D}_0 to \tilde{D}_1 parametrized by $t \in [0, 1]$ ⁴, then by the same arguments as in the previous paragraphs this leads to a family of probability distributions of the trained parameters

$$\rho_t^{par}(\theta) \quad t \in [0, 1] \quad (7)$$

²The training procedure is considered to be deterministic for simplicity

³Here and after $\mathbb{1}$ is intended to be the identity matrix and its dimension is the same as the multivariate Gaussian

⁴The continuous embedding we used is a generalization of the one presented here and is treated in Section 4.1

From now on, we will refer to t as the parameter of the continuous embedding.

We can then use the Cramér-Rao bound, stating that if we have a 1-dim family of probability distributions $\rho_t^{par}(\theta)$, smoothly parametrized by t , then the variance of any unbiased estimator \hat{t} is bounded as follows:

$$\text{var}(\hat{t}) \geq \frac{1}{I(\hat{t})} \quad (8)$$

where \tilde{t} is the real value of t and $I(\tilde{t})$ is the Fisher information defined as

$$I(\tilde{t}) = \mathbb{E}_{\rho_{\tilde{t}}^{par}(\theta)} \left[\left(\frac{\partial \log(\rho_{\tilde{t}}^{par}(\theta))}{\partial t} \Big|_{t=\tilde{t}} \right)^2 \right] \quad (9)$$

We thus have that $I(1)$ quantifies the total information available to the attacker about the presence of the TD in the training dataset, when the TD was used in the training.

The expectation value in Equation (9) involves integrating over the parameter space \mathfrak{P} , and we will see that under some assumptions, this integration leads to a sum of 1-dim integrals over the possible values of each network parameter. This allows us to consider the contribution of each network parameter to the total information, or rather, how much information about the TD is expressed in each network parameter.

4.1 Distribution of Network Parameters

To compute the Fisher information in Equation (9), we need the expression of $\rho_t^{par}(\theta)$ in a neighborhood $t \in [1 - \delta, 1]$ for some $\delta > 0$ to compute the derivative in Equation (9).

We assume that the training procedure involves the minimization of a loss function

$$\mathcal{L}(\cdot; D) : \mathfrak{P} \rightarrow \mathbb{R} \quad \text{where} \quad D \in \mathfrak{D} \quad (10)$$

so that the training function \mathcal{T} is such that for each dataset \tilde{D} the trained parameters $\tilde{\theta} = \mathcal{T}(\tilde{D})$ correspond to the minimum of $\mathcal{L}(\theta; \tilde{D})$, and thus

$$[\nabla_{\theta} \mathcal{L}](\tilde{\theta}; \tilde{D}) \equiv \frac{\partial \mathcal{L}(\theta; \tilde{D})}{\partial \theta} \Big|_{\theta=\tilde{\theta}} = 0 \quad (11)$$

and the Hessian is positive-semidefinite. We further assume that the Loss function is the average of a sample loss function over all the samples of the dataset. Thus, for any fixed dataset D and any partition of it $\{\tilde{D}_i\}_{i=1, \dots, k}$, the total loss function is

$$\mathcal{L}(\theta; D) = \sum_{i=1}^k \frac{N_i}{N} \mathcal{L}(\theta; \tilde{D}_i) \quad (12)$$

where N and N_i are the number of samples in D and in \tilde{D}_i respectively.

With these assumptions, we argue for the following radical approximation. In the limit of vanishing uncertainty ϵ of the prior of the attacker, we expect $\rho_1^{par}(\theta)$ to tend to a Dirac delta centered on $\tilde{\theta} = \mathcal{T}(\tilde{D}_1)$. Thus, for small ϵ , we expect the distribution to peak around $\tilde{\theta}$. Since the trained parameters $\tilde{\theta}$ lie in a minimum of the loss, we can approximate the loss in a small neighborhood around $\tilde{\theta}$ as

$$\mathcal{L}(\theta; \tilde{D}_1) \approx \mathcal{L}(\tilde{\theta}; \tilde{D}_1) + \frac{1}{2} (\theta - \tilde{\theta})^T \cdot H_1 \cdot (\theta - \tilde{\theta}) \quad (13)$$

where H_1 is the Hessian of the Loss for the network parameters, computed in $\tilde{\theta}$. Geometrically, this amounts to approximating the loss function with an N -dimensional paraboloid, where N is the number of network parameters. Then, intuitively, we argue that small changes to the dataset will change the loss function so that the new minimum will be displaced more along the directions where the paraboloid is wide and less along the ones where the paraboloid is narrow. More formally, we consider the fact that small perturbations of the dataset will lead to a small functional perturbation of the loss function, i.e.,

$$\mathcal{L}(\theta; \tilde{D}_1 + \delta D) \approx \mathcal{L}(\theta; \tilde{D}_1) + \delta \mathcal{L}(\theta) \quad (14)$$

Then the new minimum θ' will be displaced by a small amount according to

$$[\nabla_{\theta} \mathcal{L}](\theta'; \tilde{D}_1 + \delta D) = 0 \quad (15)$$

$$\Rightarrow H_1 \cdot (\theta' - \tilde{\theta}) + [\nabla_{\theta} \delta \mathcal{L}](\theta') = 0 \quad (16)$$

$$\Rightarrow \theta' - \tilde{\theta} = -H_1^{-1} \cdot [\nabla_{\theta} \delta \mathcal{L}](\theta') \quad (17)$$

And we see that the displacement depends on the inverse of the Hessian, so that the wider the paraboloid is along some axis, the larger the displacement will be in that direction.

Then, for a peaked Gaussian distribution around \tilde{D}_1 , we expect the distribution of the network parameters to be approximately

$$\rho_1^{par}(\boldsymbol{\theta}) \approx \mathcal{N}(\boldsymbol{\theta}; \tilde{\boldsymbol{\theta}}, C(\epsilon) H_1^{-1}) \quad (18)$$

where $C(\epsilon)$ is a scaling factor for the covariance matrix only dependent on ϵ such that it vanishes as $\epsilon \rightarrow 0$.

To compute $\rho_t^{par}(\boldsymbol{\theta})$ in the neighborhood, we first need to fix a continuous embedding for t . To simplify computations, we use a generalization of what was exposed in Section 4. We use the fact that the loss has the property of Equation (12) to carry out continuous embedding by weighting the TD on average, instead of changing the dataset. We keep fixed the distribution of the dataset ρ_1^{data} of Equation (3) so that any dataset sampled can be partitioned into D_T , sampled from ρ_T^{data} , and its complement D_{NT} (non-target), sampled from ρ_0^{data} . Then we redefine the loss function as follows

$$\mathcal{L}_t(\boldsymbol{\theta}; D) = \frac{N_{NT} \mathcal{L}(\boldsymbol{\theta}; D_{NT}) + t N_T \mathcal{L}(\boldsymbol{\theta}; D_T)}{N} \quad (19)$$

where N_{NT} and N_T are the numbers of non-target and target samples, respectively. Notice that $\mathcal{L}_1(\boldsymbol{\theta}; D) = \mathcal{L}(\boldsymbol{\theta}; D)$ and $\mathcal{L}_0(\boldsymbol{\theta}; D) = \mathcal{L}(\boldsymbol{\theta}; D_{NT})$ so that in the two extremes, one recovers the original loss if the distribution of the underlying dataset was ρ_1^{data} and ρ_0^{data} , respectively. By linearity, it follows that the Hessian of the weighted loss evaluated at $\tilde{\boldsymbol{\theta}}$ is

$$H_t = \frac{N_{NT} H(\tilde{D}_{NT}) + t N_T H(\tilde{D}_T)}{N} \quad (20)$$

Finally, in a small neighbourhood of $t = 1$, we approximate $\mathcal{T}_t(\tilde{D}) \approx \mathcal{T}(\tilde{D}) = \tilde{\boldsymbol{\theta}}$ so that by the same arguments of the previous paragraph, we extend Equation (18) to a small neighbourhood around $\tilde{\boldsymbol{\theta}}$

$$\rho_t^{par}(\boldsymbol{\theta}) \approx \mathcal{N}(\boldsymbol{\theta}; \tilde{\boldsymbol{\theta}}, C^2(\epsilon) H_t^{-1}) \quad \text{for } t \in [1 - \delta, 1] \quad (21)$$

The expression of Equation (21) requires the Hessian to be non-degenerate to be well-defined.

To allow computation, we further approximated the Hessian with its diagonal and ignored the parameters corresponding to vanishing diagonal elements, thus considering⁵

$$H_t = \text{diag}(h_1(t), \dots, h_{N_p}(t)) \quad h_i > 0 \text{ for } i = 0, \dots, N_p \quad (22)$$

where N_p is the number of network parameters after excluding the ones with vanishing diagonal elements. These approximations are used as working assumptions and are ultimately validated by experiments.

The final form of ρ_t^{par} is thus:

$$\rho_t^{par}(\boldsymbol{\theta}) \approx \prod_{i=1}^N \mathcal{N}(\theta_i; \tilde{\theta}_i, \sigma_i(t)) \quad \text{where} \quad \sigma_i(t) = \frac{C(\epsilon)}{\sqrt{h_i(t)}} \quad (23)$$

4.2 Fisher Information

Starting from the results of the previous section, we derive an explicit expression for the Fisher information of Equation (9), and define the Target Information Score.

The Fisher information about the presence of the TD, in the case where it was used during the training, is:

$$I(1) = \mathbb{E}_{\rho_1^{par}(\boldsymbol{\theta})} \left[\left(\frac{\partial \log(\rho_t^{par}(\boldsymbol{\theta}))}{\partial t} \Big|_{t=1} \right)^2 \right] \quad (24)$$

In Section 9 we show that substituting the expression of Equation (23) in the previous one we get

$$I(1) \propto \sum_{i=1}^N \left(\frac{h'_i}{h_i} \right)^2 \quad \text{with} \quad h'_i \equiv \frac{\partial h_i(t)}{\partial t} \Big|_{t=1} \quad h_i \equiv h_i(1) \quad (25)$$

where the proportionality constant only depends on ϵ . Then from Equation (20) we find:

$$h_i(t) = \frac{1}{N} \left[N_{NT} \frac{\partial^2 \mathcal{L}(\tilde{D}_{NT})}{\partial \theta_i^2} + t N_T \frac{\partial^2 \mathcal{L}(\tilde{D}_T)}{\partial \theta_i^2} \right] \quad (26)$$

So that:

$$h_i = \frac{\partial^2 \mathcal{L}(\tilde{D})}{\partial \theta_i^2} \quad h'_i = \frac{N_T}{N} \frac{\partial^2 \mathcal{L}(\tilde{D}_T)}{\partial \theta_i^2} \quad (27)$$

⁵From now on, all the derivations in the network parameters space are intended to be evaluated at the trained parameters $\tilde{\boldsymbol{\theta}}$

The total Fisher information of Equation (24) is a sum over the network parameters, and then we define the Target Information Score as

$$\text{TIS}(\theta_i) \equiv \left(\frac{\frac{\partial^2 \mathcal{L}(\tilde{D}_T)}{\partial \theta_i^2}}{\frac{\partial^2 \mathcal{L}(\tilde{D})}{\partial \theta_i^2}} \right)^2 \propto \left(\frac{h'_i}{h_i} \right)^2 \quad (28)$$

where the proportionality constant only depends on ϵ , N , N_T and is therefore independent of i .

5 Proposed Unlearning Algorithm

The framework presented in Section 4 is used to implement an FU algorithm.

We consider n clients with datasets D_β for $\beta = 1, \dots, n$, and assume that a model θ has been trained cooperatively. Suppose that the server receives an unlearning request that requires it to forget a certain *target dataset* D_T . Interestingly, we do not need D_T to be composed of samples from the same client, and instead assume it as a generic subset of the whole training dataset; thus, each client dataset will be decomposed into $D_\beta = D_\beta^{(f)} + D_\beta^{(r)}$, where $D_\beta^{(f)} = D_\beta \cap D_T$ and $D_\beta^{(r)} = D_\beta \setminus D_\beta^{(f)}$. This allows us to address the Sample Unlearning task as described in Section 3.1.

All clients are assumed to be collaborative and are required to compute the Hessian diagonal of the loss averaged on the samples $D_\beta^{(f)}$ and $D_\beta^{(r)}$, respectively, thus obtaining

$$\mathbf{h}_\beta^{(f)} = \frac{\partial^2 \mathcal{L}(D_\beta^{(f)})}{\partial \theta^2} \quad \mathbf{h}_\beta^{(r)} = \frac{\partial^2 \mathcal{L}(D_\beta^{(r)})}{\partial \theta^2} \quad \text{for } \beta = 1, \dots, n \quad (29)$$

These are then collected and aggregated by the server, allowing it to compute the expressions of Equation (27) as follows:

$$\mathbf{h} = \frac{1}{N} \sum_{\beta=1}^n \left(N_\beta^{(f)} \mathbf{h}_\beta^{(f)} + N_\beta^{(r)} \mathbf{h}_\beta^{(r)} \right) \quad \mathbf{h}' = \frac{1}{N_T} \sum_{\beta=1}^n N_\beta^{(r)} \mathbf{h}_\beta^{(f)} \quad (30)$$

where $N_\beta^{(f)}$ and $N_\beta^{(r)}$ are the number of target and non-target samples in D_β respectively. Finally, the TIS for each parameter is calculated as in Equation (9).

Then, the parameters deemed “the most informative” for D_T are identified and reset. We introduce a hyperparameter termed the *removal percentage* (α_{removal}) and in each layer, we select the α_{removal} percent of the parameters that have the highest information score. These parameters are then reset to their initial (pre-training) random values.

Finally, to recover model performance on the remaining data, a fine-tuning phase can optionally be performed. This phase employs a custom module, TRIM (Targeted Retraining via Index Masking), detailed in Algorithm 1. This module wraps the base model and defines new and different learnable parameters corresponding *only* to the scalar elements of the original parameters that were identified as informative and reset. The non-reset parameters are kept fixed (frozen) by storing them as non-trainable buffers. Retraining is performed for a single epoch, based on empirical evaluations. During this minimal federated retraining, only the reset parameters are updated using data exclusively from the remaining dataset, thus reintegrating them into the model, minimizing the influence of D_T ⁶.

5.1 Time and Space Overheads

Evaluating computational efficiency is essential for federated unlearning, as low time and space overheads are crucial for scalability and rapid response.

Time Overheads

The core of this routine involves iterating once through the complete dataset of all participating clients (i.e., processing a total of B_{total} batches across all clients). For each batch, the diagonal of the Hessian matrix (or Generalized Gauss-Newton matrix, GGN) is computed. Crucially, leveraging the BackPACK library [27] for this step, the calculation for each batch is highly efficient. The time complexity to calculate the diagonal Hessian (or GGN) for a single batch using BackPACK is approximately $k \cdot C_{bwd}$, where C_{bwd} is the cost of a standard gradient backpropagation pass for that batch, and k is a small constant factor. Therefore, this dominant step has a complexity of approximately $O(B_{\text{total}} \cdot C_{bwd})$.

Subsequent steps in the routine primarily involve $O(L \cdot N_c \cdot M_p)$ operations, where L is the number of parameter blocks (e.g., layers), N_c is the number of clients, and M_p is the cost of elementwise operations on tensors of a parameter block’s size.

⁶In fact, some influence from D_T could leak from the non-reset parameters since they also contribute to the loss landscape.

Once the information scores are computed, the unlearning process itself involves the following.

1. Identifying parameters to reset based on these precomputed scores (negligible time complexity, typically involving thresholding).
2. Reset these identified parameters (also in negligible time).
3. Optionally, a single federated retraining epoch involving only the reset parameters and the remaining clients.

This design ensures that the most computationally intensive part (Hessian diagonal computation) can be performed efficiently and on a single basis (as soon as the training ends). At the same time, the unlearning execution upon request remains ready.

Space Overheads

In terms of space overhead, the benefit of our model is related to the fact that it does not require the saving of any information during the training; nonetheless, just after the training (or as soon as the unlearning request is received), the algorithm needs the computation of the elements previously described and the saving of these.

Each client will have to compute the diagonal of the Hessian and store it until the server requires it to proceed with the unlearning. This element has approximately a space complexity equal to $O(P)$ where P is the total number of parameters.

Other elements required during training have a negligible overhead space.

6 Evaluation

6.1 Experimental setup

Dataset and models

We designed our experiments on three standard datasets: MNIST [28], FashionMNIST [29] and CIFAR10 [30]. Each dataset consists of 10 distinct classes; we consider a scenario involving five clients participating in the training process. The datasets are split into training and test sets (with a ratio of 70/30, respectively). Client-specific datasets are constructed by allocating samples from the training set according to two distinct label distribution strategies. The strategies hereby presented have been implemented following the Sample Unlearning and the Class Unlearning objectives, as outlined in Section 3:

- **Random Class Setting.** In this setting, the original dataset is randomly partitioned into 10 client subsets, corresponding to the scenario in which clients have independent and identically distributed (IID) datasets.
- **Preferential Class Setting (Categorical Unlearning Simulation).** In this setting, we model a scenario of highly heterogeneous (non-IID) distributions of the client’s datasets. For each dataset, the samples of 5 classes were distributed evenly between the five clients, while all the samples of the remaining five classes were assigned to a different client. Therefore, all clients share five classes, but each one also brings one class exclusive to that client. This setup corresponds to scenarios where individual clients hold some specialized, near-unique knowledge; in fact, unlearning any client also requires completely unlearning its preferential class.

6.1.1 Unlearning Scenario: Backdoor Attack

To evaluate the resistance of our framework against malicious participants, we simulate a backdoor attack [31]. This scenario is a practical and challenging test for data forgetting, as the unlearning algorithm must remove some deliberately inserted malicious behavior while restoring the performance.

The target client acts maliciously by inserting a backdoor trigger into its training samples. Specifically, we use the Adversarial Robustness Toolbox (ART) [32] to insert a 3x3 white pixel pattern in the corner of the images. For every sample now containing this trigger, the target client changes its label to a single predetermined target class, regardless of its original class. When these poisoned data are used to train the global model, the model learns to associate the trigger pattern with the target label. In this case, we inserted the backdoor in the MNIST dataset, making the model misclassify every such sample as the number "9", following [21]. A successful unlearning process must therefore not only neutralize the influence of this backdoor but also ensure a correct identification of the image with the correct label.

6.1.2 Training details

Regarding the models used for the experiments, we used a Convolutional Neural Network (CNN) for the tests on MNIST, composed of two convolutional layers, and we used ResNet18 ([33]) for the other two datasets.

The models were trained for 40 epochs using the ResNet18 architecture and for six epochs with CNN. Training used SGD as an optimizer, with a learning rate $\eta_{lr} = 0.01$ and a momentum of 0.9. The cross-entropy loss function was used, and a batch size of 32 was used consistently in all experiments.

For the retraining step, performed by the TRIM module, we employed the same cross-entropy loss function, and retraining of the selected parameters (and resetting) was carried out for a single epoch.

All experimental evaluations were conducted on a server equipped with an NVIDIA L40S GPU of 48 GB and an AMD 32-core CPU with 768 GB of RAM.

6.2 Evaluation Metrics

We assess the efficacy of our proposed unlearning algorithm considering three main dimensions [10]: the *efficiency* of the algorithm, the *performance* of the unlearned model, and the *forgetting* of the target dataset. The metrics used to evaluate these aspects involve the comparison of: (1) the unlearned model θ^u , (2) the model before unlearning θ^i , (3) the model retrained from scratch without the target dataset θ^r (gold standard).

Efficiency Metrics

Efficiency consists of the time and computational complexity of the unlearning process. To assess efficiency in our experiments, we use the **Recovery Time Ratio (RTR)**. This is defined as the ratio of the time required to retrain θ^r from scratch to the time taken by our unlearning algorithm:

$$\text{RTR} = \frac{\text{Retraining Time (s)}}{\text{Unlearning Time (s)}} \quad (31)$$

As discussed in Section 5, the unlearning time of our method is independent of the number of training epochs. Then we also introduce the **Break Even Epochs (BEE)**, which means the number of retrain epochs that would lead to the same computational time as our method:

$$\text{BEE} = \frac{\text{Unlearning Time (s)}}{\text{Single Epoch Retrain Time (s)}} \quad (32)$$

Given the number of epochs as E , if the required retraining epochs exceed BEE, our methods gain an efficiency advantage that grows roughly linearly with the excess $E - \text{BEE}$.

6.2.1 Performance Metrics

Model performance refers to the retained performance of the model after unlearning has occurred. In an optimal scenario, θ^u would achieve test accuracy identical to that of θ^r . The difference between the accuracy of θ^i and the accuracy of θ^r quantifies the performance impact attributable solely to the removal of the unlearned data, independent of the efficiency of the unlearning algorithm.

To further facilitate comparison, we also report a **Normalized Test Accuracy (NTA)**. This metric expresses the test accuracy of the unlearned model (θ^u) as a fraction of the benchmark retrained model's (θ^r) test accuracy:

$$\text{Normalized Test Accuracy} = \frac{\text{Accuracy}(\theta^u, D_{\text{test}})}{\text{Accuracy}(\theta^r, D_{\text{test}})} \quad (33)$$

A value approaching 1.0 for this normalized accuracy indicates that the unlearned model has successfully recovered the performance level of an ideally retrained model. NTA can also exceed 1—especially with small forget/retain sets or under non-IID FL—because the constrained update space induces pruning-like regularization, improving generalization over a full retrain.

6.2.2 Forgetting Metrics

Forgetting metrics quantify the removal of the influence of the target data (D_{target}) from the unlearned model.

The primary measure for assessing this, especially in non-IID client dataset scenarios or for categorical unlearning, is the **Target Dataset Accuracy**. We compare this accuracy across θ^u , θ^r , and θ^i . For different client datasets, we usually expect $\text{Accuracy}(\theta^i, D_T) > \text{Accuracy}(\theta^r, D_T)$, as θ^i was trained on D_T while θ^r was not. The optimal scenario is for θ^u to achieve the accuracy of the target identical to θ^r . As client datasets become more IID, the difference between

$\text{Accuracy}(\theta^i, D_T)$ and $\text{Accuracy}(\theta^r, D_T)$ may decrease, making this absolute measure less discriminative for certain settings of unlearning.

To provide a comparative measure of how effectively the unlearning process bridges the gap in knowledge between the original and the ideally retrained model, we introduce a **Normalized Forgetting Score (NFS)** based on target accuracy:

$$\text{NFS}_{\text{TargetAcc}} = 1 - \frac{|\text{Accuracy}(\theta^u, D_T) - \text{Accuracy}(\theta^r, D_T)|}{|\text{Accuracy}(\theta^i, D_T) - \text{Accuracy}(\theta^r, D_T)|} \quad (34)$$

An NFS value approaching 1.0 indicates that the unlearned model θ^u 's performance on the target data is very close to that of the ideally retrained model θ^r , relative to the initial 'vulnerability' demonstrated by θ^i . A value near zero would suggest θ^u still behaves much like θ^i on the target data.

In scenarios with IID client data, we also employ Member Inference Attacks (MIA), as described in [1]. We measure the **MIA Accuracy**, which reflects the attacker's ability to distinguish target dataset samples from test dataset samples based on model outputs. Ideally, $\text{MIA_Acc}(\theta^r) \approx 0.5$ (random guessing), while $\text{MIA_Acc}(\theta^i)$ may be significantly higher. The goal is for $\text{MIA_Acc}(\theta^u)$ to also approach 0.5. Analogous to the target accuracy, we can define a **Normalized Forgetting Score based on MIA Accuracy (NFS_{MIA})**:

$$\text{NFS}_{\text{MIA}} = 1 - \frac{|\text{MIA_Acc}(\theta^u) - \text{MIA_Acc}(\theta^r)|}{|\text{MIA_Acc}(\theta^i) - \text{MIA_Acc}(\theta^r)|} \quad (35)$$

Again, a value approaching 1.0 signifies that the unlearned model's vulnerability to MIA has been reduced to the level of the ideally retrained model, relative to the original model's vulnerability.

6.2.3 Backdoor Accuracy

For the backdoor scenario, we introduce two specific metrics to provide a complete picture of both the attack's potency and the unlearning's effectiveness:

1. **Attack Success Rate (ASR):** This metric measures the effectiveness of the backdoor attack. It is calculated as the accuracy of the model on test images containing the trigger, where the labels are set to the poisoned target label (e.g., class '9'). A high ASR in the initially trained model (θ^z) indicates a successful attack. The goal of unlearning is to reduce the ASR of the unlearned model (θ^u) to near-random guessing (e.g., 10% for a 10-class problem).
2. **Backdoor Accuracy:** We use this metric to evaluate the model's ability to correctly classify triggered examples according to their original true labels. A low Backdoor Accuracy for the θ^z model confirms that it is susceptible to the trigger. A high Backdoor Accuracy for the unlearned θ^u model demonstrates that it has successfully ignored the trigger and recovered the samples' true class, signifying successful neutralization of the backdoor.

7 Results

We experimentally assess the proposed framework by repeating each test 20 times to allow for more precise statistics. The results are presented, indicating the mean and standard deviation after the 20 tests.

Random Class (Client Unlearning)

First, we evaluate client-level unlearning, where the goal is to remove all data from a target client. Tables 1 and 2 present the performance and privacy metrics, respectively, for this scenario with 40% unlearning. As shown in Table 1, the Test Accuracy of the *Retrained* model (73.01% for CIFAR-10, 90.59% for FashionMNIST) is close to that of the *Benchmark*, reflected in NTA values of 0.94 and 0.97, respectively. Regarding privacy, Table 2 indicates that for Target Accuracy the *Retrained* model achieves NFS values of 0.99 (CIFAR-10) and 0.69 (FashionMNIST), while for MIA Accuracy, NFS_{MIA} is 0.86 (CIFAR-10) and 0.85 (FashionMNIST).

Preferential Class (Categorical Unlearning)

Next, we present results for categorical unlearning via preferential data distributions. Tables 3 and 4 summarize the results of 40% of unlearning. Table 3 shows that the *Retrained* model achieves Test Accuracies of 70.24% (CIFAR-10) and 81.59% (FashionMNIST), producing NTA values of 0.96 and 0.98. For privacy, as seen in Table 4, the NFS based on Target Accuracy is 1.00 for CIFAR-10 and 0.98 for FashionMNIST.

⁷Backdoor Accuracy is computed exclusively on triggered images whose original label was not the target class (class 9), thereby testing only the model's capacity to ignore the malicious pattern.

Table 1: Performance Metrics (40% unlearning)

Performance				
Dataset	Test Accuracy			
	Original	Benchmark	Reset	Retrained
MNIST	99.22 ± 0.07	99.13 ± 0.11	96.63 ± 2.03	99.32 ± 0.04
CIFAR10	79.48 ± 0.39	77.67 ± 0.47	10.00 ± 0.01	73.01 ± 1.27
FMNIST	93.61 ± 0.12	93.18 ± 0.22	10.03 ± 0.08	90.59 ± 0.39
NTA				
MNIST				1
CIFAR10			0.94 ± 0.02	
FMNIST			0.97 ± 0.01	

Table 2: Privacy Metrics

Privacy				
Target Accuracy				NFS
Original	Benchmark	Reset	Retrained	Retrained
99.80 ± 0.09	99.00 ± 0.15	96.98 ± 2.10	99.84 ± 0.03	1
99.95 ± 0.04	77.57 ± 0.61	10.08 ± 0.13	77.79 ± 1.65	0.99 ± 0.08
99.95 ± 0.04	77.57 ± 0.61	10.08 ± 0.13	77.79 ± 1.65	0.69 ± 0.12
MIA Accuracy				NFS _{MIA}
50.83 ± 0.32	49.84 ± 0.78	49.69 ± 0.47	49.55 ± 0.17	0.4 ± 1.27
68.48 ± 0.84	49.91 ± 0.38	50.17 ± 0.14	52.6 ± 0.35	0.86 ± 0.03
55.1 ± 0.4	49.78 ± 0.45	50.21 ± 0.12	50.58 ± 0.13	0.85 ± 0.04

Table 3: Performance Metrics (40% unlearning)

Performance				
Dataset	Test Accuracy			
	Original	Benchmark	Reset	Retrained
MNIST	99.25 ± 0.05	90.41 ± 0.09	64.25 ± 9.74	90.52 ± 0.05
CIFAR-10	79.75 ± 0.52	73.41 ± 0.36	10.91 ± 1.36	70.24 ± 0.38
FMNIST	93.58 ± 0.24	83.50 ± 0.12	10.02 ± 0.05	81.59 ± 0.65
NTA				
MNIST				1 ± 0.001
CIFAR-10				0.96 ± 0.007
FMNIST				0.98 ± 0.008

Table 4: Privacy Metrics (40% unlearning)

Privacy				
Target Accuracy				NFS
Original	Benchmark	Reset	Retrained	Retrained
99.09 ± 0.40	0.00 ± 0.00	0.24 ± 0.40	0.01 ± 0.03	1
68.28 ± 4.04	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	1
98.81 ± 0.24	0.00 ± 0.00	0.00 ± 0.00	2.41 ± 7.23	0.98 ± 0.07

Backdoor Attack Scenario

In our backdoor attack scenario (Table 5), the trained model showed complete vulnerability (100.00% ASR, 0.00% Backdoor Accuracy). After applying our unlearning method, the retrained model effectively neutralized the attack, reducing the ASR to 0.31% and recovering the Backdoor Accuracy to 98.82%. This indicates a total removal of the trigger’s influence and complete performance restoration on the original data.

Table 5: Backdoor Attack and Unlearning Metrics (40% unlearning on MNIST)

Backdoor Neutralization on MNIST				
Metric (%)	Original	Benchmark	Reset	Retrained
Attack Success Rate (ASR)	100.00 ± 0.00	0.13 ± 0.05	0.11 ± 0.17	0.31 ± 0.31
Backdoor Accuracy	0.00 ± 0.00	0.13 ± 0.05	60.16 ± 11.03	98.82 ± 0.29

Efficiency

As a final evaluation criterion, we assess time efficiency by comparing unlearning to full retraining. In Table 6, we present the results related to the efficiency metrics discussed above.

Table 6: Efficiency Metric measured on different datasets

	Efficiency	
	RTR	BEE
MNIST	0.9	7
CIFAR-10	0.34	113
FashionMNIST	0.35	117

8 Discussion

The experimental results provide insight into the capabilities of the methodology algorithm in various unlearning scenarios.

In the client unlearning scenario with random class distributions (approximating IID data; see Table 1 and Table 2), our *Retrained* model demonstrates strong performance recovery. The Normalized Test Accuracy (NTA) values, across all datasets, indicate that the performance of the unlearned model on general test data is close to that of the *benchmark* model. This suggests minimal degradation of overall model performance despite the removal of a client’s data.

We then evaluate the privacy provided by our algorithm using Normalized Forgetting Scores (NFS). For Target Accuracy, NFS values of 1 (MNIST), 0.99 (CIFAR-10) and 0.69 (FashionMNIST) indicate a substantial reduction in the *Retrained* model’s ability to correctly classify the specific ‘forgotten’ client’s data, moving its performance much closer to the *Benchmark*’s behavior on the same data. The lower NFS for FashionMNIST Target Accuracy might indicate that its more distinct features make perfect ‘as-if-never-seen’ performance harder to achieve on specific instances, even if general patterns are forgotten.

For Membership Inference Attack (MIA) Accuracy, the NFS_{MIA} values of approximately 0.85-0.86 indicate that our unlearning process significantly closes the gap between the *Original* model’s MIA vulnerability and the *Benchmark*’s near-random guessing level. This suggests effective instance-level forgetting. The initial privacy challenge in this IID-like setting, as indicated by the gap between Original and Benchmark MIA, was moderate, and our method effectively mitigates the additional risk posed by the unlearned client’s data. In MNIST, the NFS on MIA has an uncertainty of ± 1.27 , signaling that the measure is not significant in this case. This is due to the MIA Accuracy being already very similar between the original (50.83 ± 0.32) and benchmark (49.84 ± 0.78) models.

The preferential class setting (Tables 3 and 4) simulates the categorical unlearning task, where a client possesses unique knowledge. The NTA values for test accuracy (0.96 for CIFAR-10, 0.98 for FashionMNIST, and 1 for MNIST) show that even after removing a client with unique data, our *Retrained* model maintains high performance relative to the *Benchmark*.

Privacy is then assessed by the model’s inability to recognize the forgotten unique class. The Target Accuracy for the *Retrained* model drops to near zero, as evidenced by the NFS values of 1.00 (CIFAR-10 and MNIST) and 0.98 (FashionMNIST). This signifies almost perfect forgetting and indicates negligible residual predictive power in the forgotten class.

To validate our framework’s resilience against targeted security threats, we simulated a poisoning-based backdoor attack. The results offer clear evidence of the defensive capabilities of our method. The initially trained model was compromised entirely, exhibiting a perfect Attack Success Rate (ASR) of 100.00% while its Backdoor Accuracy on the true labels dropped to 0.00%. This confirms that the model had learned to prioritize the malicious trigger above all else, making it untrustworthy. However, the unlearning process was able to perform a complete reversal. The intermediate reset model demonstrates that it is sufficient to destroy the backdoor, reducing the ASR to a negligible 0.11% while already recovering 60.16% of the original classification ability. However, a greater increase has been observed following the implementation of the retrained model, where, after only a single fine-tuning epoch, the model not only kept the ASR near-zero (0.31%) but also restored the Backdoor Accuracy to a near-perfect 98.82%. Collectively, these findings demonstrate that our unlearning framework serves as a potent and efficient defense mechanism, capable of precisely removing malicious behavior from a federated model and fully restoring its integrity.

Regarding efficiency, the Recovery Time Ratio (RTR) indicates that computational times for our unlearning method are higher than the full retraining in all experiments. This represents a limitation of our model that has been highlighted in our tests due to the small number of epochs needed to train the base models—when retraining is fast, time savings from unlearning are modest. However, our method scales better: as tasks get harder and models require more epochs to converge, retraining cost grows with E, while our unlearning cost grows much more slowly, meaning that the efficiency gap widens. BEE provides a valuable metric to gauge when our unlearning approach is projected to reach a break-even

point in terms of temporal cost compared to full retraining. For ResNet18, the measured BEE is 113 in CIFAR10 and 117 on FashionMNIST.

8.1 Sensitivity Analysis

As the final step, we investigated the effect of varying the percentage of *removal* (α_{removal}), which dictates the extent of the parameter reset as defined in Section 5. Our evaluations focused on the resulting model metrics as previously described.

The experiment presented here was conducted on the CIFAR-10 dataset, following the setup as in Section 6.1. The MIA was evaluated using the logistic model.

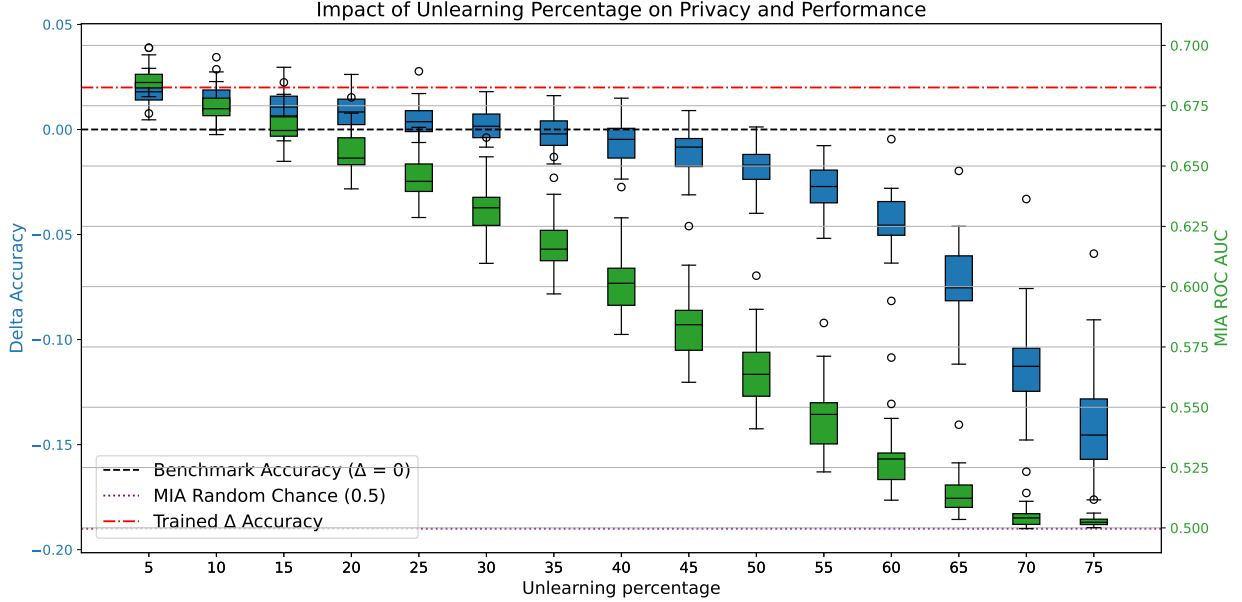


Figure 1: Effect of Unlearning Percentage (α_{removal}). Left y-axis (blue): Delta test accuracy of the unlearned model relative to the benchmark; right y-axis (green): MIA accuracy of the unlearned model.

Increasing α_{removal} , thereby resetting (and subsequently retraining) a larger portion of influential parameters, generally led to a lower MIA accuracy on the ‘forgotten’ data. At high α_{removal} values, the accuracy of MIA approached random chance (approximately 0.5), signifying strong erasure of instance-level data. We observe that higher values of α_{removal} are generally associated with a reduction in the predictive accuracy of the model on the general test set. These findings, illustrated in Figure 1, highlight how the degree of parameter reset influences both the level of data forgetting and the overall utility of the model.

9 Conclusion

In this paper, we have addressed the challenge of efficient and effective data erasure in federated learning by introducing a principled, information-theoretic framework for federated unlearning. Our approach models the deletion of a target dataset as a parameter-estimation problem; by leveraging second-order statistics, we can quantify the contribution of each parameter to data memorization and selectively reset only the most influential parameters before performing a single-epoch lightweight fine-tuning operation in the remaining clients. This targeted strategy yields privacy outcomes and performance levels that closely approximate those of complete retraining, with an efficiency that is constant in the number of training epochs.

Our empirical evaluations demonstrate the practical balance achieved by this two-stage process. The initial *Reset* of identified parameters leads to a significant degradation in the performance metrics related to the forgotten data, confirming the neutralization of influential parameters. The subsequent minimal retraining phase effectively recovers the utility of the global model while essentially preserving the privacy gains achieved. The percentage of unlearning

naturally serves as a key hyperparameter to navigate this balance based on specific application needs, as we have shown in Section 8.1.

Furthermore, our framework exhibits several desirable properties for real-world deployment: it is both data- and architecture-agnostic (the algorithm works with arbitrary datasets), and significantly, it operates without the need to store historical model updates or extensive checkpoint snapshots. These characteristics make it a well-suited and scalable solution for implementing "right-to-be-forgotten" mandates in federated systems.

However, we argue that a significant limitation of our framework, which must be addressed for real-world deployment, is its efficiency. In the tests conducted so far, this has been a problematic aspect of our model. A second limitation is our reliance on a diagonal approximation for the Hessian computation, where we ignore parameters corresponding to vanishing diagonal elements.

Future Work

Future research will focus on key extensions to our framework and on addressing its current limitations.

From an efficiency standpoint, we see two possibilities. The first is a technical exploration into more optimized methods for Hessian computation. The second is to conduct experiments on larger models trained on tasks that require more epochs, which would permit us to verify the hypothesis of increased efficiency due to scaling factors.

Other future work could include exploring the utility of richer, potentially more precise, second-order information beyond Hessian diagonals (e.g., block-diagonal or approximations of the full Hessian). We also aim to investigate adaptive hyperparameter tuning for the unlearning percentages and to analyze the resulting privacy-performance trade-offs.

Appendix

Appendix 1

We start by computing

$$\frac{\partial}{\partial t} \log(\rho_t^{par}(\boldsymbol{\theta})) = \sum_{i=1}^N \frac{\partial}{\partial t} \log(\mathcal{N}(\theta_i; \tilde{\theta}_i, \sigma_i(t))) \quad (36)$$

$$= \sum_{i=1}^N \left[\frac{\partial}{\partial \sigma_i} \log(\mathcal{N}(\theta_i; \tilde{\theta}_i, \sigma_i(t))) \right] \frac{\partial \sigma_i}{\partial t} \quad (37)$$

A simple computation shows that

$$\frac{\partial}{\partial \sigma_i} \log(\mathcal{N}(\theta_i; \tilde{\theta}_i, \sigma_i(t))) \quad (38)$$

$$= \frac{1}{\mathcal{N}(\theta_i; \tilde{\theta}_i, \sigma_i(t))} \frac{\partial}{\partial \sigma_i} \left(\frac{1}{\sqrt{2\pi} \sigma_i(t)} \exp \left[\frac{(\theta_i - \tilde{\theta}_i)^2}{2\sigma_i(t)^2} \right] \right) \quad (39)$$

$$= -\frac{1}{\sigma_i(t)} + \frac{(\theta_i - \tilde{\theta}_i)^2}{\sigma_i(t)^3} \quad (40)$$

and

$$\frac{\partial \sigma_i}{\partial t} = -\frac{1}{2} C(\epsilon) (h_i(t))^{-\frac{3}{2}} \frac{\partial h_i(t)}{\partial t} \quad (41)$$

$$= -\frac{1}{2} \frac{\sigma_i(t)^3}{C(\epsilon)^2} \frac{\partial h_i(t)}{\partial t} \quad (42)$$

so that we find

$$\frac{\partial}{\partial t} \log(\rho_t^{par}(\boldsymbol{\theta})) = \frac{1}{2C(\epsilon)^2} \sum_{i=1}^N (\sigma_i(t)^2 - (\theta_i - \tilde{\theta}_i)^2) \frac{\partial h_i(t)}{\partial t} \quad (43)$$

Then we get

$$\left(\left. \frac{\partial \log(\rho_t^{par}(\boldsymbol{\theta}))}{\partial t} \right|_{t=1} \right)^2 \quad (44)$$

$$= \frac{1}{4C^4(\epsilon)} \left\{ \sum_{i \neq j} h'_i h'_j [\sigma_i(1)^2 \sigma_j(1)^2 - \sigma_i(1)^2 (\theta_j - \tilde{\theta}_j)^2 - (\theta_i - \tilde{\theta}_i)^2 \sigma_j(1)^2 + (\theta_i - \tilde{\theta}_i)^2 (\theta_j - \tilde{\theta}_j)^2] + \sum_{i=1}^N [\sigma_i(1)^4 - 2(\theta_i - \tilde{\theta}_i)^2 \sigma_i(1)^2 + (\theta_i - \tilde{\theta}_i)^4] \right\} \quad (45)$$

where we wrote $h'_i \equiv \left. \frac{\partial h_i(t)}{\partial t} \right|_{t=1}$. We can then compute the expectation value

$$I(1) = \mathbb{E}_{\rho_1^{par}(\boldsymbol{\theta})} \left[\left(\left. \frac{\partial \log(\rho_t^{par}(\boldsymbol{\theta}))}{\partial t} \right|_{t=1} \right)^2 \right] \quad (46)$$

$$= \frac{1}{2C^4(\epsilon)} \sum_{i=1}^N \sigma_i(1)^4 h_i'^2 \quad (47)$$

$$= \frac{1}{2C^4(\epsilon)} \sum_{i=1}^N \left(\frac{h'_i}{h_i} \right)^2 \quad (48)$$

where we wrote $h_i \equiv h_i(1)$ and we used the gaussian moments

$$\mathbb{E}_{\rho_1^{par}(\boldsymbol{\theta})} [(\theta_i - \tilde{\theta}_i)^2] = \sigma_i(1)^2 \quad (49)$$

$$\mathbb{E}_{\rho_1^{par}(\boldsymbol{\theta})} [(\theta_i - \tilde{\theta}_i)^4] = 3\sigma_i(1)^4 \quad (50)$$

Appendix 2

Appendix 2 describes the algorithm implemented for the retraining: 1 and 2

Algorithm 1 Efficient Targeted Retraining using the TRIM wrapper.

```

1: Module State:
2:  $\mathcal{F}$  {The underlying (frozen) model architecture/logic.}
3: BUFFERS {Map: ParamName  $\rightarrow$  Tensor; stores the non-trainable parameter base values.}
4: TRAINABLE_PARAMS {ParameterDict: ParamName  $\rightarrow$  nn.Parameter; holds the trainable elements.}
5: INDICES {Map: ParamName  $\rightarrow$  Tensor; stores indices specifying trainable element locations.}
6: function INITIALIZE( $\Theta_{\text{base}}, M_{\text{info}}$ )
7:   Input:
8:      $\Theta_{\text{base}}$  {Parameters of the original, trained base model.}
9:      $M_{\text{info}}$  {Map {ParamName  $\rightarrow$  Tensor of indices for retraining}.}
10:
11:   Output: Initialized BUFFERS, TRAINABLE_PARAMS, INDICES
12:
13:   Store underlying model architecture  $\mathcal{F}$ .
14:   BUFFERS  $\leftarrow \{\}$ 
15:   TRAINABLE_PARAMS  $\leftarrow \{\}$ 
16:   INDICES  $\leftarrow \{\}$ 
17:
18:   for each (name,  $\theta$ ) in  $\Theta_{\text{base}}$  do
19:      $\theta_{\text{fixed}} \leftarrow \theta.\text{clone}()$  {Copy original parameter values into a “fixed” buffer.}
20:     if name  $\in M_{\text{info}}$  then
21:       indices  $\leftarrow M_{\text{info}}[\text{name}]$ 
22:        $k \leftarrow \text{length}(\text{indices})$ 
23:       if  $k > 0$  then
24:          $\theta_{\text{train}} \leftarrow$  New ‘nn.Parameter’ of size  $k$  {Only these  $k$  entries will be trainable.}
25:         Initialize  $\theta_{\text{train}}$  (e.g. with zeros or original values).
26:         TRAINABLE_PARAMS[name]  $\leftarrow \theta_{\text{train}}$ 
27:         INDICES[name]  $\leftarrow$  indices
28:          $\theta_{\text{fixed}}[\text{indices}] \leftarrow 0$  {Zero-out base values where training will occur.}
29:       end if
30:     end if
31:     Register  $\theta_{\text{fixed}}$  into BUFFERS[name] {Store as a non-trainable buffer.}
32:   end for
33: end function

```

Algorithm 2 Forward pass for the TRIM wrapper.

```

1: function FORWARD( $X$ )
2:   Input:  $X$  {Input data batch.}
3:   Output: Model output using reconstructed parameters.
4:
5:    $\Theta_{\text{live}} \leftarrow \{\}$  {Construct current parameters for this forward pass.}
6:
7:   for each parameter name ‘name’ stored in BUFFERS do
8:      $\theta_p \leftarrow$  BUFFERS[name].clone() {Start with the fixed (zeroed-out) base parameter.}
9:     if name  $\in$  TRAINABLE_PARAMS then
10:       Retrieve indices  $\leftarrow$  INDICES[name]
11:       Retrieve  $\theta_{\text{train}} \leftarrow$  TRAINABLE_PARAMS[name]
12:        $\theta_p[\text{indices}] \leftarrow \theta_p[\text{indices}] + \theta_{\text{train}}$  {Add trainable values into their correct positions.}
13:     end if
14:      $\Theta_{\text{live}}[\text{name}] \leftarrow \theta_p$ 
15:   end for
16:    $Y_{\text{out}} \leftarrow \mathcal{F}(X; \text{params} = \Theta_{\text{live}})$  {Run the underlying architecture using reconstructed live parameters.}
17:   return  $Y_{\text{out}}$ 
18: end function

```

Acknowledgments

Computational resources provided by hpc-ReGAINs@DISCo, which is a MUR Department of Excellence Project within the Department of Informatics, Systems and Communication at the University of Milan-Bicocca (<https://www.disco.unimib.it>).

Author contributions

A.B. and L.M. designed the model and the computational framework and analysed the data. A.B. and L.M. carried out the implementation. A.B. and L.M. wrote the manuscript with input from all authors. F.S. and A.M. provided critical feedback and were in charge of overall direction and planning. All authors discussed the results and commented on the manuscript.

Data availability

Code for this paper is available at: <https://github.com/lorenzomanini/FedUnlearn-PE>

References

- [1] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models, 2017.
- [2] Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej Kos, and Dawn Song. The secret sharer: Evaluating and testing unintended memorization in neural networks, 2019.
- [3] Paul Voigt and Axel Bussche. *The EU General Data Protection Regulation (GDPR): A Practical Guide*. Springer Cham, 01 2017.
- [4] Rik Crutzen, Gjalt-Jorn Ygram Peters, and Christopher Mondschein. Why and how we should care about the general data protection regulation. *Psychology & Health*, 34(11):1347–1357, 2019. PMID: 31111730.
- [5] California Civil Code. Assembly bill no. 1008, 2023.
- [6] Preston Bukaty. *The California Privacy Rights Act (CPRA) – An implementation and compliance guide*. IT Governance Publishing, 2021.
- [7] H. Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. Communication-efficient learning of deep networks from decentralized data, 2023.
- [8] Lucas Bourtoutle, Varun Chandrasekaran, Christopher A. Choquette-Choo, Hengrui Jia, Adelin Travers, Baiwu Zhang, David Lie, and Nicolas Papernot. Machine unlearning. In *2021 IEEE Symposium on Security and Privacy (SP)*, pages 141–159, 2021.
- [9] Gaoyang Liu, Xiaoqiang Ma, Yang Yang, Chen Wang, and Jiangchuan Liu. Federaser: Enabling efficient client-level data removal from federated learning models. In *2021 IEEE/ACM 29th International Symposium on Quality of Service (IWQOS)*, pages 1–10, 2021.
- [10] Nicolò Romandini, Alessio Mora, Carlo Mazzocca, Rebecca Montanari, and Paolo Bellavista. Federated unlearning: A survey on methods, design guidelines, and evaluation metrics. *IEEE Transactions on Neural Networks and Learning Systems*, page 1–21, 2024.
- [11] Lucas Bourtoutle, Varun Chandrasekaran, Christopher A Choquette-Choo, Hengrui Jia, Adelin Travers, Baiwu Zhang, David Lie, and Nicolas Papernot. Machine unlearning. In *2021 IEEE symposium on security and privacy (SP)*, pages 141–159. IEEE, 2021.
- [12] Shangchao Su, Bin Li, and Xiangyang Xue. One-shot federated learning without server-side training. *Neural Networks*, 164:203–215, July 2023.
- [13] Yinzhi Cao and Junfeng Yang. Towards making systems forget with machine unlearning. In *2015 IEEE symposium on security and privacy*, pages 463–480. IEEE, 2015.
- [14] Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In *International conference on machine learning*, pages 1885–1894. PMLR, 2017.
- [15] Chuan Guo, Tom Goldstein, Awni Hannun, and Laurens Van Der Maaten. Certified data removal from machine learning models. *arXiv preprint arXiv:1911.03030*, 2019.
- [16] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR, 2017.
- [17] Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Advances and open problems in federated learning. *Foundations and trends® in machine learning*, 14(1–2):1–210, 2021.

- [18] Chen Wu, Sencun Zhu, and Prasenjit Mitra. Federated unlearning with knowledge distillation. *arXiv preprint arXiv:2201.09441*, 2022.
- [19] Guanhua Ye, Tong Chen, Quoc Viet Hung Nguyen, and Hongzhi Yin. Heterogeneous decentralised machine unlearning with seed model distillation. *CAAI Transactions on Intelligence Technology*, 9(3):608–619, 2024.
- [20] Matthew Jagielski, Milad Nasr, Christopher Choquette-Choo, Katherine Lee, and Nicholas Carlini. Students parrot their teachers: Membership inference on model distillation, 2023.
- [21] Anisa Halimi, Swanand Kadhe, Ambrish Rawat, and Nathalie Baracaldo. Federated unlearning: How to efficiently erase a client in fl?, 2023.
- [22] Vitaly Feldman. Does learning require memorization? A short tale about a long tail. *CoRR*, abs/1906.05271, 2019.
- [23] Aditya Golatkar, Alessandro Achille, and Stefano Soatto. Eternal sunshine of the spotless net: Selective forgetting in deep networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9304–9312, 2020.
- [24] Zhuo Ma, Yang Liu, Ximeng Liu, Jian Liu, Jianfeng Ma, and Kui Ren. Learn to forget: Machine unlearning via neuron masking. *IEEE Transactions on Dependable and Secure Computing*, 20(4):3194–3207, 2022.
- [25] Davis W. Blalock, Jose Javier Gonzalez Ortiz, Jonathan Frankle, and John V. Guttag. What is the state of neural network pruning? *CoRR*, abs/2003.03033, 2020.
- [26] Anvith Thudi, Gabriel Deza, Varun Chandrasekaran, and Nicolas Papernot. Unrolling sgd: Understanding factors influencing machine unlearning. In *2022 IEEE 7th European Symposium on Security and Privacy (EuroS&P)*, pages 303–319. IEEE, 2022.
- [27] Felix Dangel, Frederik Kunstner, and Philipp Hennig. Backpack: Packing more into backprop. *CoRR*, abs/1912.10985, 2019.
- [28] Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012.
- [29] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *CoRR*, abs/1708.07747, 2017.
- [30] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical Report 0, University of Toronto, Toronto, Ontario, 2009.
- [31] Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Identifying vulnerabilities in the machine learning model supply chain, 2019.
- [32] Maria-Irina Nicolae, Mathieu Sinn, Minh Ngoc Tran, Beat Buesser, Ambrish Rawat, Martin Wistuba, Valentina Zantedeschi, Nathalie Baracaldo, Bryant Chen, Heiko Ludwig, Ian M. Molloy, and Ben Edwards. Adversarial robustness toolbox v1.0.0, 2019.
- [33] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CoRR*, abs/1512.03385, 2015.