

Efficient Unlearning with Privacy Guarantees

Josep Domingo-Ferrer, Najeeb Moharram Jebreel, and David Sánchez

Universitat Rovira i Virgili, Department of Computer Engineering and Mathematics,
CYBERCAT - Center for Cybersecurity Research of Catalonia,
Av. Països Catalans 26, 43007 Tarragona, Catalonia
{josep.domingo, najeeb.jebreel, david.sanchez}@urv.cat

July 8, 2025

Abstract

Privacy protection laws, such as the GDPR, grant individuals the right to request the forgetting of their personal data not only from databases but also from machine learning (ML) models trained on them. Machine unlearning has emerged as a practical means to facilitate model forgetting of data instances seen during training. Although some existing machine unlearning methods guarantee exact forgetting, they are typically costly in computational terms. On the other hand, more affordable methods do not offer forgetting guarantees and are applicable only to specific ML models. In this paper, we present *efficient unlearning with privacy guarantees* (EUPG), a novel machine unlearning framework that offers formal privacy guarantees to individuals whose data are being unlearned. EUPG involves pre-training ML models on data protected using privacy models, and it enables *efficient unlearning with the privacy guarantees offered by the privacy models in use*. Through empirical evaluation on four heterogeneous data sets protected with k -anonymity and ϵ -differential privacy as privacy models, our approach demonstrates utility and forgetting effectiveness comparable to those of exact unlearning methods, while significantly reducing computational and storage costs. Our code is available at <https://github.com/najeebjebreel/EUPG>.

Keywords: Machine unlearning, Privacy, Differential privacy, k -Anonymity, Right to be forgotten.

1 Introduction

Personal data are commonly used to train machine learning (ML) models. However, due to privacy concerns, individuals can request the removal of their data used to train a model. The legal basis for these requests is the right to be forgotten, which is guaranteed by several data protection regulations, such as the European Union’s General Data Protection Regulation (GDPR [21]), the California Consumer Privacy Act (CCPA [45]), and the recently approved European AI Act [36]. In particular, this right remains valid even if the user initially consented to the use of their data for ML.

State-of-the-art privacy attacks against trained ML models have shown that it is not only possible to discover whether a particular item was part of the training data (through membership inference attacks [10,42]), but the training data themselves can be reconstructed (partially or totally) from the model using reconstruction attacks [2,23,49]. Hence, forgetting one item in the training data requires modifying the models trained on those data.

A naive strategy to achieve data forgetting consists in retraining the model from scratch, after excluding from the training data the items requested for removal. However, doing this is computationally expensive.

Machine unlearning (MU) [5, 9] aims to provide a cheaper alternative to retraining from scratch, while maintaining the utility of the model. Depending on the influence of the unlearned data items on the unlearned model, MU methods can be classified as exact and approximate [4, 35, 53]. The exact unlearning methods [5, 9, 54] aim to completely eliminate the influence of the data elements that must be unlearned, ensuring that the model behaves as if it had never seen those elements during training. Although this approach offers formal forgetting guarantees, it typically incurs high computational and storage costs, sometimes close to the cost of retraining from scratch. On the other hand, approximate unlearning methods [3, 25, 28, 41, 46] approximate complete removal of target items through efficient strategies to update the parameters of the ML model. However, these methods lack formal forgetting guarantees and only provide approximate guarantees specific to certain ML model types, which makes them insufficient for strict compliance of the right to be forgotten. A promising middle ground is DP-certified unlearning [13, 15], which probabilistically limits data influence using (ϵ, δ) -DP guarantees. However, existing DP-certified approaches suffer from restrictive assumptions (*e.g.*, strong convexity), utility degradation in non-convex settings, and scalability limitations.

In this paper, we propose a novel framework for machine unlearning named *efficient unlearning with privacy guarantees* (EUPG), which provides formal privacy guarantees to individuals whose data are unlearned. Our framework modifies both the training data and the initial model training to enable efficient and utility-preserving unlearning. As a result, it significantly reduces computational and storage costs in comparison with exact unlearning methods, while still providing privacy guarantees derived from the privacy models enforced on the training data. Moreover, our approach accommodates a wider range of ML model types than most current unlearning techniques.

In summary, our contributions are:

- An efficient machine unlearning framework that guarantees the privacy of individuals whose data are being unlearned. Our framework can be combined with *any* privacy model.
- A first instantiation of our framework with k -anonymity [38] as a privacy model. The resulting method guarantees that the unlearned data are protected by k -anonymity.
- A second instantiation with differential privacy (DP) [17] as a privacy model. The resulting method guarantees that the unlearned data are protected by ϵ -DP.
- Experiments and detailed analyses on four data sets, involving tabular data and images, that demonstrate that our framework achieves utility and forgetting effectiveness comparable to those of exact unlearning methods, while incurring significantly lower computational and storage unlearning costs.

The remainder of this paper is organized as follows. Section 2 discusses related work on machine unlearning. Section 3 provides background on the privacy models that we use in our instantiations. Section 4 presents our general framework for efficient unlearning with privacy guarantees. Section 5 reports empirical results on a variety of data sets and ML models. Section 6 draws conclusions and sketches some lines for future research.

2 Related work

Machine unlearning attempts to produce an ML model having unlearned some data item ideally as if it had been retrained from scratch without that item, but at a lower cost.

Exact unlearning aims to completely eliminate the influence of data items to be forgotten by using strategies that are more efficient than retraining from scratch. SISA [5] was proposed as a generic framework for exact unlearning through algorithmic-level retraining. This framework divides the data set into S shards of R slices and trains partial ML models slice-by-slice with checkpoints. Upon completion, partial models are

saved within their respective shards. During inference, the outputs are aggregated in an ensemble style. For data forgetting, the slice containing the data to be forgotten is identified, the target data are removed, and the ML model is retrained from the last checkpoint to ensure effective unlearning. SISA ensures exact forgetting and is also adaptable to a wide range of ML tasks and model architectures. However, its practicality is limited by the large computational and memory costs associated with saving models and checkpoints, re-training, and inference, which become even more acute in the case of large language models [4]. Moreover, the utility might be compromised because of the smaller training data employed to build the independent models on disjoint shards. Various improvements and adaptations of SISA for different models have been proposed, including k -means [24], extremely randomized trees [40], random forests [6,40], and graph-based models [11]. Nevertheless, each adaptation is specifically designed for a particular model type. Machine unlearning through algorithmic stability [48] leverages total variation (TV) stability and coupling techniques to achieve exact unlearning, primarily for convex risk minimization problems. While this approach can be extended to non-convex models like deep neural networks (DNNs), its guarantees are weaker in such settings, and practical implementation may require additional heuristics, such as gradient clipping.

Approximate unlearning aims to provide a practical alternative to exact unlearning by efficiently minimizing the influence of unlearned items to an acceptable level through an efficient and *ex post* alteration of the learned model [53]. Many approximate unlearning methods do not offer formal guarantees of data removal, and the evidence of removal is at best empirical [28, 44, 47, 50, 56].

DP-certified unlearning could also be viewed as approximate unlearning, as it probabilistically limits rather than completely eliminates the influence of the data to be forgotten. However, it offers certified privacy through (ϵ, δ) -DP guarantees. It is gradient-based, relying on noisy optimization during training for unlearning. [15] use noisy gradient descent to fine-tune the trained model on the remaining data after excluding the data to be forgotten, ensuring privacy by making the unlearned model indistinguishable from the one never trained on the deleted data. However, noisy fine-tuning may degrade model utility, especially in non-convex settings. [13] formalize Rényi unlearning under strong convexity, achieving exponential decay in privacy loss with projected noisy stochastic gradient descent. However, their approach assumes strong convexity, which limits applicability to non-convex settings. [14] introduce Langevin unlearning, unifying DP learning and unlearning via noisy gradients. While they optimize privacy-utility-complexity trade-offs, particularly in convex regimes, the theoretical bounds for non-convex problems are not yet tight, limiting the applicability of their approach in such scenarios. More similar to our approach, [32] propose Rewind-to-Delete (R2D), a first-order black-box method designed for non-convex functions, which achieves (ϵ, δ) -certified unlearning through a three-step process: i) rewinding to an intermediate model checkpoint saved during the original training process, ii) retraining the model from the checkpoint for a small number of iterations on the retained data set (*i.e.*, the data set excluding the data to be unlearned), and iii) adding Gaussian noise to the model weights to make the unlearned model statistically indistinguishable from one retrained from scratch on the retained data. This approach circumvents convexity assumptions and is thus applicable to a wide range of ML models, including DNNs; it also maintains computational efficiency by avoiding full retraining. However, R2D requires the loss function to be Lipschitz-smooth, which may not hold for all models or data sets. In addition, its privacy guarantees weaken as more data points are unlearned, potentially requiring more noise to maintain privacy.

Our proposed framework introduces a novel approach by modifying the training data and the pipeline to enable unlearning that is more efficient than retraining or SISA-based methods. Unlike approximate unlearning methods, it offers formal guarantees of data removal. In contrast to DP-certified methods, which are limited to approximate (ϵ, δ) -DP guarantees and place requirements on the ML model, EUPG is agnostic both to the ML model and to the privacy guarantee: it can be applied to any ML model and with any privacy model (privacy guarantee). For example, it can offer any form of DP (including pure ϵ -DP, (ϵ, δ) -DP, Rényi DP, etc.), k -anonymity or any k -anonymity extensions that protect not only against re-identification but also against attribute disclosure (such as l -diversity, t -closeness, etc.). Moreover, in contrast to gradient-centric

DP-certified methods, our framework is plug-and-play via data modification, which decouples unlearning from model architecture. Once data are transformed via the selected privacy guarantee, any ML model (gradient-based or not) trained on the transformed data inherits the privacy guarantee for the unlearned data.

3 Background

This section provides background on data anonymization and the main two families of privacy models (*i.e.*, privacy guarantees).

Tabular data sets, also called microdata sets, are often used to train ML models. A tabular data set is composed of records, where each record reports several attributes on an individual. Considering their disclosure potential, attributes can play the following roles: (i) *identifiers* (*e.g.*, passport no., social security no., name-surname, etc.) unequivocally identify the individual to whom the record corresponds; (ii) *quasi-identifiers* (*e.g.*, zipcode, profession, gender, age, etc.) do not identify the individual when taken separately, but they may when considered jointly (*e.g.*, a 90-year old female doctor living in a certain zipcode is probably unique); (iii) *confidential attributes* (*e.g.*, diagnosis, income, etc.) contain sensitive information; iv) other attributes may exist that are not in the previous three groups.

When data are anonymized, the identifiers are suppressed. However, this is not enough, and some modification of the quasi-identifiers is needed to prevent *re-identification* by an adversary. To offer *ex ante* guarantees against re-identification, this modification should be done under the scope of a *privacy model*.

k-Anonymity [38] was the first privacy model proposed, and it was followed by several variants and extensions (like *l*-diversity [30] or *t*-closeness [29], that additionally protect against disclosure of confidential attributes). A data set is *k*-anonymous if, for each combination of quasi-identifiers present in the data set, there are at least *k* records sharing that combination. These records form a so-called *k*-anonymous class. In this way, the probability of successful re-identification is at most $1/k$.

An original data set can be made *k*-anonymous in several ways, including:

- *Generalization and suppression*. This was the method originally proposed by the inventors of *k*-anonymity [38]. For categorical quasi-identifiers, categories are coarsened (*e.g.*, a profession “Medical doctor” can be coarsened into “Healthcare personnel”, a zipcode “08001” can be coarsened into “0800*”). Numerical quasi-identifiers can be coarsened into intervals (*e.g.*, age 33 into a 30-34 age interval). Outlier values can be suppressed to reduce the need for too harsh a coarsening.
- *Microaggregation* [16]. This method microaggregates records by their quasi-identifier attributes, that is, creates clusters of at least *k* records (and less than $2k$ records) such that the quasi-identifier values within each cluster are maximally similar. Then, the centroid of the quasi-identifier values in each cluster is computed and used to replace the quasi-identifier values of the records in the cluster. In the Anatomy or probabilistic *k*-anonymity variant [43, 52], the computation of centroids is avoided; instead, the quasi-identifier combinations within each cluster are randomly permuted.

ϵ -Differential privacy (DP, [17]) is a privacy model that inaugurated another family consisting of several variants or relaxations [20]. DP seeks to bound the influence of any particular record in the original data set (and thus of the individual to whom the record corresponds) on the statistical outcomes. Formally speaking, a randomized query function κ satisfies ϵ -differential privacy if, for all data sets D_1 and D_2 that differ in one record (also known as neighbor data sets), and all $S \subset \text{Range}(\kappa)$, we have $\Pr(\kappa(D_1) \in S) \leq \exp(\epsilon) \Pr(\kappa(D_2) \in S)$.

For a numerical query f , ϵ -DP can be reached via noise addition, that is, by computing $\kappa(x) = f(x) + N$, where N stands for noise. Typically, the Laplace noise distribution is used. The amount of noise that needs to be added is inversely proportional to ϵ (smaller ϵ means greater privacy and thus requires more noise) and

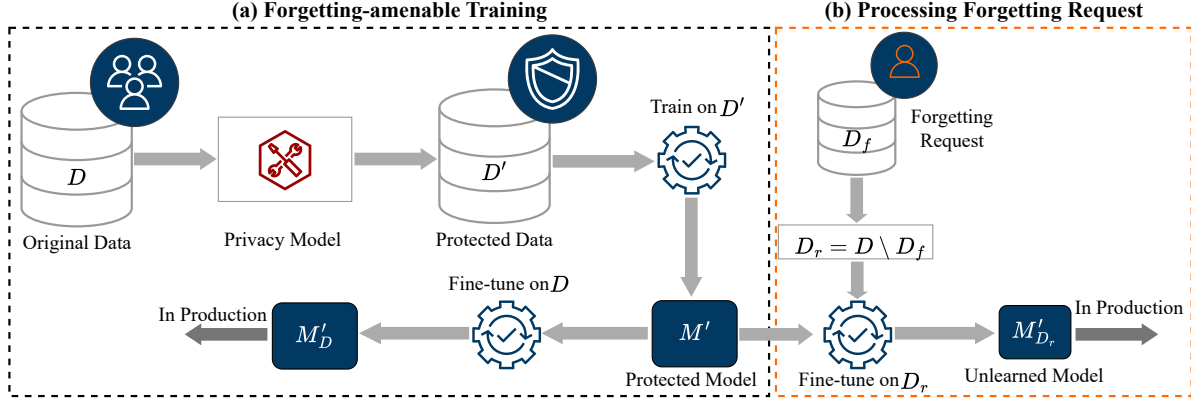


Figure 1: Workflow of the proposed EUPG framework

directly proportional to the sensitivity of the query function f (variability between neighbor data sets), that is, $\Delta_f = \max_{D_1, D_2 \in \mathcal{D}} \|f(D_1) - f(D_2)\|_1$, where D_1 and D_2 are two data sets differing in one record and \mathcal{D} is the collection of data sets on which f can be evaluated.

For non-numerical queries, the exponential mechanism [31] is typically used. This mechanism selects outputs with a probability proportional to the exponent of the utility of the outputs, ensuring that more “useful” results are more likely, while still protecting privacy.

DP has some interesting properties:

- *Immunity to post-processing.* If a function f provides ϵ -DP, then any function g , applied to the output of f , also preserves ϵ -DP.
- *Sequential composition.* Let κ_1 be a randomized function that satisfies ϵ_1 -DP and κ_2 a randomized function that satisfies ϵ_2 -DP. Then, any deterministic function of (κ_1, κ_2) satisfies $(\epsilon_1 + \epsilon_2)$ -DP.
- *Parallel composition.* Let κ_1 and κ_2 be randomized functions that satisfy ϵ -DP. If κ_1 and κ_2 are applied to *disjoint* data sets, any deterministic function of (κ_1, κ_2) satisfies ϵ -DP.

Although DP was initially proposed to protect queries to a remote database, it can also be applied to anonymize tabular microdata sets by enforcing the Laplacian mechanism or the exponential mechanism on *all* the attributes of the microdata set [39, 44]. In this case, the global sensitivity is the maximum variability of each attribute (typically, the attribute range). Moreover, since attribute values within a record are typically correlated, sequential composition applies: hence, the budget ϵ needs to be split among the attributes in the data set, which means that the noise added to each attribute increases with the number of attributes.

4 Efficient unlearning with privacy guarantees

Our framework for *efficient unlearning with privacy guarantees* (EUPG) is designed to simultaneously address the efficiency-privacy-utility aspects of unlearning. Using privacy models and engineering of the training process, it enables effective and efficient subsequent unlearning with privacy guarantees. The general workflow of EUPG is shown in Figure 1 and it consists of two main stages: (a) forgetting-amenable training and (b) processing of forgetting requests.

Forgetting-amenable training This stage involves preparing the ML model in a way that inherently supports efficient and effective unlearning. First, a privacy model offering formal privacy guarantees is chosen

(such as those introduced in Section 3 or a variant/extension of them). The choice of the privacy model depends on multiple factors, including the data types of the attributes and the specific privacy guarantees that are sought. For a particular privacy model, stronger privacy guarantees need stricter parameters (*e.g.*, larger k for k -anonymity or smaller ϵ for DP). Once a privacy model is chosen, the original data set (\mathbf{D}) is transformed into a protected version (\mathbf{D}') by a mechanism that enforces the privacy model (as described in Section 3). An initial ML model M' is then pre-trained on \mathbf{D}' , which transfers the formal privacy guarantees to M' . Note that immunity to post-processing is not necessarily invoked here: since the guarantees hold against an attacker being given full access to \mathbf{D}' , in particular, they hold if the attacker chooses to train a model M' on \mathbf{D}' . This initial data protection and pre-training steps are the most computationally expensive part of our framework, but they are one-time processes that set the stage for efficient, utility-preserving unlearning with privacy guarantees. Then M' undergoes fine-tuning on the entire original data set \mathbf{D} to mitigate the potential degradation in the utility of the ML model due to privacy-preserving modifications to the training data. Fine-tuning allows the protected ML model to quickly make up for the utility lost during anonymization (which caused M' to learn only general patterns without memorizing item-level information). The resulting fine-tuned model $M'_{\mathbf{D}}$ is then deployed for production.

Processing of forgetting requests Upon receiving a forgetting request, EUPG can handle it efficiently without retraining the ML model. First, the ML model currently deployed is deleted and the model M' pre-trained on the data set \mathbf{D}' protected by the privacy model is fine-tuned on the retained data $\mathbf{D}_r = \mathbf{D} \setminus \mathbf{D}_f$, where \mathbf{D}_f are the data to be forgotten. The resulting fine-tuned model $M'_{\mathbf{D}_r}$ is then deployed for production. Since elements in \mathbf{D}_f have not individually contributed to the fine-tuned model ($M'_{\mathbf{D}_r}$), the privacy of the individual(s) whose data are in \mathbf{D}_f is protected by design under the guarantees of the privacy model embedded in the pre-trained protected model (M'). In addition, the use of fine-tuning significantly reduces the computational and storage costs with respect to those typically required for retraining in exact unlearning. The fine-tuning step triggered by a forgetting request is key in our approach, as it provides the model manager with the flexibility to balance model utility against the computational cost of unlearning: if the manager prioritizes utility (resp. low cost), they can increase (resp. decrease) the number of fine-tuning epochs.

In the following, we show how this general framework can be instantiated with the two families of privacy models introduced in Section 3, and we prove the resulting privacy guarantees.

4.1 Unlearning with a k -anonymity privacy guarantee

Under k -anonymity, the probability that an adversary can re-identify a forgotten record after unlearning based on the quasi-identifiers is upper bounded by $1/k$.

Algorithm 1 (k -Anonymity amenable training).

1. Let \mathbf{D} be a data set used for training. Let QI be the set of quasi-identifier attributes in \mathbf{D} .
2. k -Anonymize \mathbf{D} into \mathbf{D}^k by using the most convenient computational method depending on the nature of the quasi-identifier attributes QI (see Section 3).
3. Train a machine learning model on \mathbf{D}^k . Let the trained model be M^k .
4. Fine-tune M^k on \mathbf{D} to obtain a model $M^k_{\mathbf{D}}$.

Protocol 1 (k -Anonymous unlearning).

1. An individual asks the manager of model $M^k_{\mathbf{D}}$ to remove her data \mathbf{D}_f from \mathbf{D} and $M^k_{\mathbf{D}}$.

2. The model manager deletes $M_{\mathbf{D}}^k$ and fine-tunes M^k on $\mathbf{D}_r = \mathbf{D} \setminus \mathbf{D}_f$ to obtain $M_{\mathbf{D}_r}^k$.

Proposition 1 (Privacy). *Unlearning with Protocol 1 satisfies k -anonymity.*

Proof: Protocol 1 falls back on M^k , which has been trained on the k -anonymous version \mathbf{D}^k of the original data set \mathbf{D} . Since k -anonymity holds against an attacker \mathcal{A} with full access to \mathbf{D}^k (i.e., the attacker’s probability of successful re-identification via quasi-identifiers is at most $1/k$), in particular, it holds if \mathcal{A} decides to train M^k on \mathbf{D}^k and hence M^k does not break k -anonymity (no immunity to post-processing property needs to be invoked).

Then, the fine-tuning to obtain $M_{\mathbf{D}_r}^k$ takes place on all items of \mathbf{D} except those in the forget set \mathbf{D}_f . Therefore, the information in \mathbf{D}_f is only retained by $M_{\mathbf{D}_r}^k$ through its k -anonymous class in \mathbf{D}^k , and thus k -anonymous privacy is maintained for the individual who requested \mathbf{D}_f to be forgotten. \square

k -Anonymity can be replaced by any of its extensions that protect against confidential attribute disclosure (such as l -diversity or t -closeness), and it is straightforward to adapt Algorithm 1, Protocol 1 and Proposition 1.

4.2 Unlearning with a differentially private guarantee

We next show how to obtain an ϵ -DP guarantee for the forgetting request, which means that after unlearning the forgotten data should be unnoticeable from the output of the ML model except by a factor $\exp(\epsilon)$.

Algorithm 2 (DP-amenable training).

1. Let \mathbf{D} be a data set used for training.
2. Transform \mathbf{D} into an ϵ -DP data set \mathbf{D}^ϵ (e.g., by applying the Laplacian or exponential mechanisms to every attribute; see Section 3). If record attributes are not independent and thus sequential composition applies, use a budget $\epsilon/|\text{attributes}|$ for each attribute to obtain ϵ -DP protected records.
3. Train a machine learning model on \mathbf{D}^ϵ . Let the trained model be M^ϵ .
4. Fine-tune M^ϵ on \mathbf{D} to obtain a model $M_{\mathbf{D}}^\epsilon$.

Protocol 2 (ϵ -DP unlearning).

1. An individual asks the manager of model $M_{\mathbf{D}}^\epsilon$ to remove her data \mathbf{D}_f from \mathbf{D} and $M_{\mathbf{D}}^\epsilon$.
2. The model manager fine-tunes M^ϵ on $\mathbf{D}_r = \mathbf{D} \setminus \mathbf{D}_f$ to obtain $M_{\mathbf{D}_r}^\epsilon$.

Proposition 2 (Privacy). *Unlearning with Protocol 2 satisfies ϵ -DP.*

Proof: Protocol 2 falls back on M^ϵ , which has been trained on the ϵ -DP version \mathbf{D}^ϵ of the original data set \mathbf{D} . Then, the fine-tuning to obtain $M_{\mathbf{D}_r}^\epsilon$ excludes the forget set \mathbf{D}_f . According to the post-processing immunity property of DP (see Section 3), the ϵ -DP guarantee achieved by \mathbf{D}^ϵ extends to $M_{\mathbf{D}_r}^\epsilon$ for the individual(s) to whom \mathbf{D}_f corresponds. \square

ϵ -DP can be replaced by any of its variants or relaxations, and it is straightforward to adapt Algorithm 2, Protocol 2, and Proposition 2. Note that post-processing immunity holds for all state-of-the-art variants and relaxations of DP. However, the above proof could be rewritten without invoking the post-processing immunity property, but instead using an attacker argument as in the proof of Proposition 1.

5 Experimental results

We conducted experiments on four data sets and three ML models to assess the performance of our framework in terms of the preservation of utility of the ML model, the effectiveness of forgetting and computational efficiency. We also compared against retraining the model from scratch and against SISA, a method also providing formal unlearning guarantees. Comparison against approximate unlearning methods would be unfair to our framework as they do not offer guarantees of forgetting. Furthermore, we investigated how the privacy model chosen influences the performance of our approach and examined its sensitivity to various hyperparameters. The experiments were carried out on a system running Windows 11 Home OS, equipped with a 12th Gen Intel®Core™i7-12700 12-core CPU, 32 GB RAM, and an NVIDIA GeForce RTX 4080 with 16 GB GPU. Our code is available at <https://github.com/najeebjebrael/EUPG>

5.1 Experimental setup

Data sets We used four publicly available data sets, each representing a classification problem from a different domain. Three of these data sets are tabular and were chosen because of their privacy relevance, as they contain records describing personal data of individuals. In addition, we included an image classification data set to evaluate the generality of our approach across various data types.

- *Adult income*¹: It comprises 32,561 training records and 16,281 testing records of demographic and financial data, with six numerical and eight categorical attributes. The class attribute indicates whether an individual makes more than 50K dollars a year.
- *Heart disease*²: It contains 55,869 training records and 14,131 testing records of patient data, with five numerical and six categorical measurements related to cardiovascular diseases. The class attribute denotes the presence of heart disease.
- *Credit information*³: It includes 96,215 training records and 24,054 testing records of financial information, with ten numerical attributes. The class attribute indicates whether an individual has experienced financial distress.
- *CIFAR-10*⁴: It is a widely used image classification data set containing 60,000 32x32 pixel images across ten classes, with three RGB channels. The data set is divided into 50,000 training samples and 10,000 testing samples.

ML models For each tabular data set, we built two ML classification models: a multi-layer perceptron (MLP) and an XGBoost classifier. We implemented MLP models using PyTorch with an input layer, a hidden layer, and an output layer. Regarding the XGBoost classifier, we utilized the implementation provided by XGBoost [12], which can be found on the official XGBoost website⁵. For CIFAR-10, we used the DenseNet12 [27] deep model, which is known for its densely connected layers that enhance feature reuse and improve gradient flow.

¹<https://archive.ics.uci.edu/ml/datasets/Adult>

²<https://www.kaggle.com/sulianova/cardiovascular-disease-dataset>

³<https://www.kaggle.com/c/GiveMeSomeCredit>

⁴<https://www.cs.toronto.edu/~kriz/cifar.html>

⁵<https://xgboost.ai/>

Evaluation metrics The utility of the ML models was evaluated using accuracy (Acc) for the Adult, Heart and CIFAR-10 data sets. For the Credit data set, which has a significant class imbalance (only 6.92% of the data belong to the positive class), the area under the ROC curve (AUC) was used instead. The effectiveness of forgetting was assessed using the loss-based [55] and the entropy-based [37] membership inference attacks (MIAs) implemented in TensorFlow Privacy⁶. We employed AUC to evaluate MIAs [26, 51], as our objective was to quantify the trade-off between true and false positives of member and non-member data points.

Privacy models We used the two privacy models described in Section 4 for the tabular data sets. To enable a fair comparison, we applied both privacy models to all attributes except the class attribute. To implement k -anonymity, we used the MDAV microaggregation algorithm [16] with one-hot encoded categorical attributes and the parameter $k \in \{3, 5, 10, 20, 80\}$.

For DP, we applied the Laplace mechanism [19] for numerical attributes and the exponential mechanism [31] for categorical attributes, with $\epsilon \in \{0.5, 2.5, 5, 25, 50, 100\}$. For DP on the Adult data set, which contains several non-ordinal categorical attributes, we leveraged the unsupervised training capabilities of TabNet [1] to generate embeddings for these attributes, which allowed us to encode each non-ordinal categorical attribute into a 10-dimensional embedding vector. Then, we used the cosine similarity between the embeddings of attributes as a utility function for the exponential mechanism.

For the CIFAR-10 data set, we enforced DP via the DP-Pix methodology described in [22]. DP-Pix first pixelizes the image by averaging pixel values in blocks of size $b \times b$ to reduce sensitivity and then adds noise to the pixels using the Laplace mechanism based on the global sensitivity $255m/b^2$. The parameter m defines the number of different pixels between neighboring images. We used $m = 16$, as suggested by the DP-Pix author, and $b = 4$, which is consistent with image resolution [22]. We did not enforce k -anonymity on CIFAR-10 images because aggregating or generalizing images to make them indistinguishable (as required by k -anonymity) produces meaningless images for any safe (large) enough k , and very few works employing k -anonymity in images use $k = 2$ [8, 33, 34].

Training settings For each tabular data set, we trained an MLP model and an XGBoost model from scratch on the entire training set \mathbf{D} . We used cross-entropy loss and the Adam optimizer to train all MLP models. The MLP hidden layer was configured with 128 neurons for all benchmarks, with the exception of the Credit benchmark, which was configured with 256 neurons. For CIFAR-10, we trained the DenseNet12 model using cross-entropy loss and the SGD optimizer. The specific hyperparameters used during the training of these benchmarks are detailed in Table 1.

For SISA, we split the training set into 5 disjoint shards (each containing 10 slices) and applied the SISA training procedure using the training hyperparameters presented in Table 1.

For our EUPG with tabular data, we obtained the base private models M^k and M^ϵ by training them on the \mathbf{D}^k and \mathbf{D}^ϵ protected data sets, respectively, using the same training hyperparameters in Table 1. For EUPG on CIFAR-10, we trained on \mathbf{D}^ϵ only. We then fine-tuned M^k and M^ϵ on \mathbf{D} for some epochs (see details below) to obtain $M_{\mathbf{D}}^k$ and $M_{\mathbf{D}}^\epsilon$ with the learning rates and batch size presented in Table 1.

Forgetting settings We randomly sampled a forget set \mathbf{D}_f from the original training set \mathbf{D} using a given forgetting ratio (*e.g.*, 5%). After subtracting \mathbf{D}_f from \mathbf{D} , the remaining training points made up the retain set \mathbf{D}_r . To forget \mathbf{D}_f with privacy guarantees, we fine-tuned M^k and M^ϵ on \mathbf{D}_r to obtain $M_{\mathbf{D}_r}^k$ and $M_{\mathbf{D}_r}^\epsilon$, respectively. We used the learning rates and batch sizes in Table 1. For MLP and DenseNet12, we fine-tuned the models for a number of epochs. For XGBoost, we added a number of estimators and fine-tuned the models with those added estimators. See details on fine-tuning below.

⁶https://www.tensorflow.org/responsible_ai/privacy/guide

Table 1: Training hyperparameters

data set	Model	Hyperparameters
Adult income	MLP	BS:512, LR:1e-2, Epochs:100
	XGBoost	Estimators:300, Depth:10, LR:0.5, λ :5
Heart disease	MLP	BS:512, LR:1e-2, Epochs:200
	XGBoost	Estimators:200, Depth:7, LR:0.5, λ :5
Credit	MLP	BS:256, LR:1e-3, Epochs:200
	XGBoost	Estimators:200, Depth:9, LR:0.5, λ :5
CIFAR-10	DenseNet12	BS:64, Epochs:100 LR:(1e-1 for train., 1e-2 for finetun.)

For all training, forgetting, and attack experiments, we report the average results of three runs.

5.2 Main results

We present our findings on the performance of EUPG, SISA, and retraining from scratch before and after forgetting.

Performance before forgetting Tables 2 and 3 show the pre-forgetting performance of the original models trained on \mathbf{D} from scratch ($M_{\mathbf{D}}^O$), those obtained by SISA ($M_{\mathbf{D}}^{SISA}$), and those obtained by EUPG ($M_{\mathbf{D}}^k$ and $M_{\mathbf{D}}^{\epsilon}$). The former table contains results for the MLP and XGBoost models on tabular data, and the latter table contains results for the DenseNet12 models on image data. For a fair comparison with SISA, we report the results for EUPG with ‘safe’ privacy parameters (*i.e.*, $\epsilon \leq 1$, as discussed in [18], and $k = 10$, as suggested in [7]). Our experiments on the influence of hyperparameters (Section 5.3) give further empirical justification for our choices of k and ϵ .

We first discuss our results on tabular data. Our method with $k = 10$ achieved the best utility score in 3 out of 6 benchmarks and it was the second-best in the other 3 benchmarks. With $\epsilon = 0.5$, our method was the best or second-best in the MLP benchmarks, but caused utility degradation in the XGBoost benchmarks, especially for the Credit data set. This degradation comes from the sequential training and fine-tuning mechanisms of XGBoost itself, where the current tree is affected by all its preceding trees. Additionally, the XGBoost split criteria can be highly sensitive to changes in data distribution. The noise added for DP can alter the apparent distribution of the data, thereby leading to suboptimal splits that would not have been chosen in the absence of noise. Later in the fine-tuning step, the leaf values are updated by \mathbf{D}_r while the splits are fixed. All of this hampers the recovery of the model’s utility after fine-tuning. On the other hand, since k -anonymity preserves the general distribution of data, it benefits more from fine-tuning than DP.

For CIFAR-10 with DenseNet12, EUPG achieved a similar accuracy to SISA, but both underperformed compared to the original model. The fact that our method achieved higher utility scores than the original and SISA models with tabular data can be attributed to the fundamental principle underlying our approach, which is that learning from anonymized (*i.e.*, more general, less individual-specific) data mitigates the usual problem of overfitting to the training data set. Therefore, fine-tuning the base protected ML models on \mathbf{D} for a few epochs (5 in our experiments) was enough to regain utility without affecting the model’s ability to generalize to unseen data. Although this was not the case for image data with DenseNet12, we show in Section 5.3 that fine-tuning DenseNet12 for a larger number of epochs can recover all utility. The impact of the number of fine-tuning epochs on model utility is examined and discussed in Section 5.3.

Table 2: Performance of MLP and XGBoost models on tabular data before forgetting. RT denotes runtime in seconds. EUPG was fine-tuned on \mathbf{D} for 5 epochs. Best scores are in bold. Second-best are underlined.

data set	Model Method/Metric	MLP			XGBoost		
		Utility↑ Acc(%) / AUC(%)	MIA↓ AUC(%)	RT (s)↓	Utility↑ Acc(%) / AUC(%)	MIA↓ AUC(%)	RT (s)↓
Adult income	$M_{\mathbf{D}}^O$	84.22	55.13	95.51	<u>85.88</u>	54.68	4.16
	$M_{\mathbf{D}}^{SISA}$	83.91	63.16	<u>96.76</u>	<u>85.80</u>	66.60	<u>15.99</u>
	$M_{\mathbf{D}}^{k=10}$	85.53	<u>51.33</u>	149.58	86.68	51.71	56.14
	$M_{\mathbf{D}}^{\epsilon=0.5}$	<u>84.97</u>	50.76	287.17	81.57	<u>52.20</u>	195.54
Heart disease	$M_{\mathbf{D}}^O$	72.35	52.90	198.27	72.78	58.20	1.95
	$M_{\mathbf{D}}^{SISA}$	71.95	64.83	<u>201.69</u>	73.05	72.95	<u>8.51</u>
	$M_{\mathbf{D}}^{k=10}$	73.82	50.81	227.06	<u>72.85</u>	<u>51.58</u>	31.40
	$M_{\mathbf{D}}^{\epsilon=0.5}$	<u>73.52</u>	<u>50.86</u>	227.39	<u>68.92</u>	50.43	31.33
Credit	$M_{\mathbf{D}}^O$	75.71	52.66	475.36	80.80	55.72	2.34
	$M_{\mathbf{D}}^{SISA}$	75.39	72.15	361.87	83.22	70.46	<u>11.98</u>
	$M_{\mathbf{D}}^{k=10}$	<u>81.55</u>	49.99	555.43	<u>82.20</u>	<u>50.42</u>	107.12
	$M_{\mathbf{D}}^{\epsilon=0.5}$	82.09	<u>50.08</u>	496.00	56.84	50.18	47.54

Regarding privacy protection, EUPG significantly improved over the original and SISA models. This is indicated by lower MIA AUC values across all data sets and models. The reason is our method’s effectiveness in reducing the memorization of individual data points, the key factor contributing to vulnerability against MIA. We can see that SISA achieved the worst protection against MIA because the single models in SISA are trained on shards that are small compared to the whole training data (training on smaller data is known to increase overfitting). It can also be seen that the MIA AUC is often greater with XGBoost, which suggests that XGBoost models are prone to memorizing training samples than MLP models.

Regarding the runtime required to obtain the models before forgetting, SISA generally had runtimes similar to the original models with MLP and DenseNet12 benchmarks. With XGBoost benchmarks, SISA incurred higher runtimes compared to the original models, although it ranked second. Looking at the runtimes of EUPG, we can see that they were higher than those of the original and SISA models. This increase in runtime came mainly from the computational cost required to obtain \mathbf{D}^k and \mathbf{D}^ϵ , and also from the fine-tuning step. The runtime increase on Adult with $\epsilon = 0.5$ w.r.t. the other methods on Adult is notable: it is due to the added embedding time of the categorical values in this data set. Detailed runtimes of our method with different values of k and ϵ are given in Section 5.4.

In summary, compared to the original and SISA models, EUPG provides the best privacy protection while preserving the utility of models before forgetting. In particular, EUPG is more effective in mitigating MIAs before forgetting. This sometimes comes at a high computational cost, but *this cost is one time and tolerable given the greater privacy protection that our method offers* compared to the original models and SISA. We argue that privacy needs to be protected not only when forgetting is requested but also beforehand.

Performance after forgetting Tables 4 and 5 show the post-forgetting performance of the original ML models retrained on \mathbf{D}_r from scratch ($M_{\mathbf{D}_r}^O$), those obtained with SISA ($M_{\mathbf{D}_r}^{SISA}$), and those obtained with EUPG ($M_{\mathbf{D}_r}^k$ and $M_{\mathbf{D}_r}^\epsilon$) when forgetting 5% of the training data across the four data sets. The former table contains results for MLP and XGBoost models on tabular data and the latter table results for DenseNet12 models on image data.

Regarding utility, with MLP benchmarks, EUPG with $k = 10$ achieved the best utility across all data sets. With XGboost, our method with $k = 10$ achieved the best or second-best utility, but it degraded utility

Table 3: Performance of DenseNet12 models on image data before forgetting. RT denotes runtime in seconds. EUPG was fine-tuned on \mathbf{D} for 5 epochs. Best scores are in bold. Second-best are underlined.

data set	Model Method/Metric	DenseNet12		
		Utility↑ Acc(%) / AUC(%)	MIA↓ AUC(%)	RT (s)↓
CIFAR-10	$M_{\mathbf{D}}^O$	88.60	<u>63.50</u>	<u>28704.32</u>
	$M_{\mathbf{D}}^{SISA}$	<u>82.83</u>	<u>74.46</u>	27337.17
	$M_{\mathbf{D}}^{k=10}$	82.09	56.78	30737.36
	$M_{\mathbf{D}}^{\epsilon=0.5}$			

Table 4: Performance of MLP and XGBoost models on tabular data after forgetting 5% of the training data. EUPG was fine-tuned on \mathbf{D}_r for 5 epochs. Best scores are in bold. Second-best are underlined.

data set	Model Method/Metric	MLP			XGBoost		
		Utility↑ Acc(%) / AUC(%)	MIA↓ AUC(%)	RT (s)↓	Utility↑ Acc(%) / AUC(%)	MIA↓ AUC(%)	RT (s)↓
Adult income	$M_{\mathbf{D}_r}^O$	83.31	50.12	89.91	85.54	<u>50.77</u>	<u>4.09</u>
	$M_{\mathbf{D}_r}^{SISA}$	84.18	51.12	48.47	<u>85.78</u>	<u>51.64</u>	8.41
	$M_{\mathbf{D}_r}^{k=10}$	85.40	<u>50.64</u>	<u>4.65</u>	86.51	50.26	0.39
	$M_{\mathbf{D}_r}^{\epsilon=0.5}$	<u>84.92</u>	51.21	4.21	81.79	51.04	0.39
Heart disease	$M_{\mathbf{D}_r}^O$	72.77	50.12	186.73	72.30	50.73	<u>2.23</u>
	$M_{\mathbf{D}_r}^{SISA}$	71.76	50.59	96.54	73.12	50.07	4.69
	$M_{\mathbf{D}_r}^{k=10}$	73.87	51.10	4.66	<u>72.99</u>	51.03	0.24
	$M_{\mathbf{D}_r}^{\epsilon=0.5}$	<u>73.73</u>	<u>50.36</u>	<u>4.86</u>	69.03	<u>50.37</u>	<u>0.25</u>
Credit	$M_{\mathbf{D}_r}^O$	71.96	50.74	448.61	80.24	49.90	3.68
	$M_{\mathbf{D}_r}^{SISA}$	76.12	50.30	179.81	83.33	50.37	7.19
	$M_{\mathbf{D}_r}^{k=10}$	81.78	<u>50.69</u>	10.56	<u>82.24</u>	51.21	<u>0.33</u>
	$M_{\mathbf{D}_r}^{\epsilon=0.5}$	<u>81.66</u>	50.81	<u>11.45</u>	60.71	<u>49.97</u>	0.32

with $\epsilon = 0.5$. SISA and our method with $\epsilon = 0.5$ also reduced utility with CIFAR-10-DenseNet12 compared to the original model.

Regarding forgetting, all ML models achieved comparable MIA AUC scores (close to random guessing), which is consistent with the exact/guaranteed unlearning they achieve by design.

Regarding the forgetting runtime, EUPG achieved a significant improvement compared to the original and SISA models. This is clearly visible in the very short runtimes observed for both the MLP and XGBoost benchmarks.

These results underscore the core strength of our approach, namely its ability to facilitate efficient and effective forgetting while preserving the utility of the ML model. This trifecta of benefits —efficiency, effectiveness, and utility preservation— is attained by initiating the process with privacy-enhanced pre-trained models (M^k and M^ϵ). These models capture general patterns within the training data without memorizing individual data points. This allows fine-tuning those models on the retain set \mathbf{D}_r to expedite utility recovery and improvement while retaining the privacy guarantees embedded within those pre-trained models.

Storage requirements and inference efficiency SISA requires storage proportional to the product of data shards (S) and slices within each shard (R). This significantly scales up storage needs. In contrast, EUPG only stores two versions of the model: M^k and $M_{\mathbf{D}}^k$ (resp. M^ϵ and $M_{\mathbf{D}}^\epsilon$ if using DP) before forgetting, and M^k and $M_{\mathbf{D}_r}^k$ (resp. M^ϵ and $M_{\mathbf{D}_r}^\epsilon$ if using DP) afterward, resulting in lightweight and stable storage

Table 5: Performance of DenseNet12 models on image data after forgetting 5% of the training data. EUPG was fine-tuned on \mathbf{D}_r for 5 epochs. Best scores are in bold. Second-best are underlined.

data set	Model Method/Metric	DenseNet12		
		Utility↑ Acc(%) / AUC(%)	MIA↓ AUC(%)	RT (s)↓
CIFAR-10	$M_{\mathbf{D}_r}^O$	88.18	<u>51.10</u>	25833.89
	$M_{\mathbf{D}_r}^{SISA}$	<u>81.72</u>	49.66	<u>13399.18</u>
	$M_{\mathbf{D}_r}^{\epsilon=0.5}$	81.10	<u>50.73</u>	1568.50

requirements. In terms of inference efficiency, SISA incurs a computational overhead of approximately $(S - 1)T$, where T denotes the inference time per model; this degrades response times as the number of shards increases. In contrast, our approach does not introduce additional latency over the original model’s inference time.

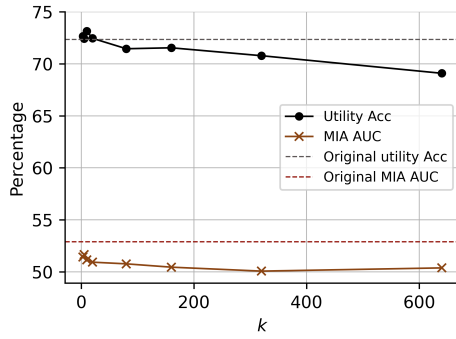
5.3 Impact of hyperparameters

In this section, we discuss the influence of key hyperparameters. For each of them, we report results on the data sets and models that best illustrate the effect of varying its value.

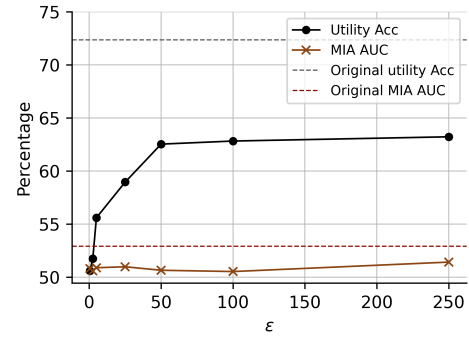
Impact of parameter k Figure 2a depicts the impact of parameter k (when using k -anonymity to protect training data) on utility and forgetting. The figure reports the performance of the MLP model trained on the anonymized Heart disease data set (without fine-tuning, M^k) for various values of k in terms of accuracy (utility) and MIA AUC scores (forgetting). The results show a noticeable trend, where the accuracy peaks at $k = 10$ and then gradually decreases as k increases. This suggests that a moderate level of anonymity (with $k = 10$) strikes an optimal balance for utility due to effective generalization of the training data without losing important features. At lower values of k ($k = 3$), the precision is slightly lower than at $k = 10$, which could be due to the insufficient generalization of the ML model. As k increases beyond 10, the utility gradually decreases. This is indicative of overgeneralization, where the ML model loses useful information to make accurate predictions. The MIA AUC score is slightly higher with $k = 3$ and 5, but the improvement in forgetting is not linearly correlated with the increase of k . However, in general, forgetting improves as k increases.

In summary, the value of k defines a trade-off between the utility of the model and the privacy (or forgetting) guarantees, in line with the typical use of k -anonymity to protect microdata releases. Therefore, it is crucial to choose a value of k that maximizes the utility of the model while providing enough privacy. Given the results in Figure 2a, we took $k = 10$ in the rest of our experiments.

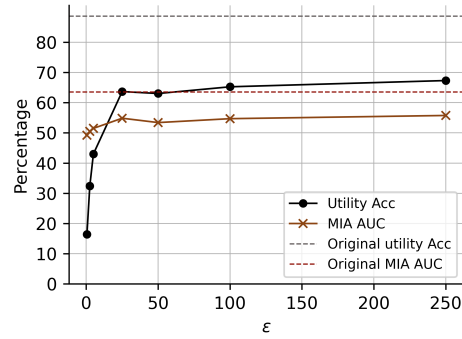
Impact of parameter ϵ Figure 2b shows the impact of the differential privacy parameter ϵ on utility and forgetting. The figure outlines the performance of MLP trained on the differentially private Heart data set (without fine-tuning, M^ϵ) for various values of ϵ in terms of accuracy (utility) and MIA AUC scores (forgetting). As ϵ increases, which corresponds to a relaxation of privacy guarantees, there is a notable improvement in the accuracy of the model. This improvement is especially significant as ϵ goes from 0.5 to 25, suggesting that moderate relaxation of privacy can substantially enhance the ability of the ML model to learn from data. However, the rate of accuracy improvement tapers off for very high values of ϵ , approaching a plateau that hints at a diminishing return on utility gains. The MIA AUC scores remain relatively stable



(a) Impact of k for the MLP model trained on the anonymized Heart data set



(b) Impact of ϵ for the MLP model trained on the anonymized Heart data set



(c) Impact of ϵ for the DenseNet12 model trained on the anonymized CIFAR-10 data set

Figure 2: Impact of privacy parameters on utility and forgetting

Table 6: Impact of the number of fine-tuning epochs on MLP and XGBoost trained on Adult with $\epsilon = 0.5$

data set	Epochs	MLP			XGBoost		
		Utility	Acc(%) \uparrow	MIA AUC(%) \downarrow	Utility	Acc(%) \uparrow	MIA AUC(%) \downarrow
Adult	0		59.89	51.42		73.25	50.97
	5		84.92	51.21		81.79	51.04
	10		84.97	51.32		82.34	51.14
	20		84.93	51.36		82.58	51.33

Table 7: Impact of the number of fine-tuning epochs on DenseNet12 trained on CIFAR-10 with $\epsilon = 0.5$

data set	Epochs	DenseNet12	
		Utility	Acc(%) \uparrow MIA AUC(%) \downarrow
CIFAR-10	0	16.42	49.23
	5	81.10	50.73
	10	85.72	50.58
	20	88.03	51.07

across the spectrum of ϵ values, indicating that even a light perturbation of the data (corresponding to a weak DP guarantee) can significantly reduce the vulnerability of the ML model to MIAs.

Figure 2c shows a similar trend for CIFAR-10 with DenseNet12. From Figure 2b and Figure 2c alone, it would seem that choosing, say, $\epsilon = 100$ would yield a better trade-off between utility and privacy than the value $\epsilon = 0.5$ we have taken in our experiments. The problem is that, although MIAs are the standard way to evaluate forgetting in the literature, they are rather weak attacks from the privacy point of view. In order to provide a fair comparison with the exact forgetting of SISA or retraining from scratch, which achieve top privacy (equivalent to taking $\epsilon = 0$ in DP), we could not take ϵ too large. We therefore chose $\epsilon < 1$ following the guidelines of [18], specifically $\epsilon = 0.5$. This level of privacy is comparable to exact forgetting and can protect against other attacks, such as reconstruction.

The above issue arises from the fact that the ϵ -DP guarantee on insensitivity to individual contributions becomes meaningless when ϵ is large. Notice that this does not happen with k -anonymity, whose privacy guarantee is meaningful for any $k \geq 2$: the re-identification probability is upper bounded by $1/k$, and thus one can choose any value of k such that $1/k$ is considered to be an acceptable re-identification probability.

Impact of the number of fine-tuning epochs Table 6 reports the impact of the number of fine-tuning epochs on utility and privacy for the MLP and XGBoost models trained on Adult with $\epsilon = 0.5$. We can see that fine-tuning significantly boosts the utility (accuracy) of both models, resulting in rapid gains in performance within the first 5 epochs. For MLP, utility jumps from 59.89% to approximately 84.92% and stabilizes with minimal changes beyond 5 epochs. For XGBoost, starting from a higher baseline of 73.25%, utility improves steadily up to 82.58% at 20 epochs.

Table 7 reports the same results for the DenseNet12 model trained on CIFAR-10. Here, the impact of fine-tuning epochs is more notable. Starting from a low baseline of 16.42% accuracy, the model shows a significant improvement, reaching 81.10% accuracy after only 5 epochs. This trend continues, with utility reaching 85.72% at 10 epochs and 88.03% at 20 epochs. For all benchmarks, MIA AUC scores remain similar to those of M^ϵ (the protected model before fine-tuning), indicating almost no impact of fine-tuning epochs on forgetting performance.

Table 8: Impact of the forgetting ratio with the Credit data set

Method	Forgetting ratio	MLP		XGBoost			
		Utility	AUC(%) \uparrow	RT (s) \downarrow	Utility	AUC(%) \uparrow	RT (s) \downarrow
$M_{\mathbf{D}_r}^O$	5%		71.96	448.61	80.24		3.68
	10%		71.46	402.03	80.12		3.60
	20%		70.92	383.03	79.76		3.52
	50%		67.11	236.01	79.55		3.25
$M_{\mathbf{D}_r}^{SISA}$	5%		76.12	179.81	83.33		7.19
	10%		76.20	160.09	81.37		7.60
	20%		76.52	146.06	80.91		7.59
	50%		76.20	94.38	79.81		7.24
$M_{\mathbf{D}_r}^{k=10}$	5%		81.78	10.56	82.24		0.33
	10%		81.87	10.32	82.01		0.30
	20%		81.95	9.73	81.97		0.26
	50%		81.61	6.17	81.74		0.19
$M_{\mathbf{D}_r}^{\epsilon=0.5}$	5%		81.66	11.45	60.71		0.32
	10%		81.89	10.32	60.34		0.29
	20%		81.89	9.53	56.84		0.26
	50%		81.92	6.21	52.52		0.18

Impact of the forgetting ratio Table 8 reports the impact of the forgetting ratio on utility and runtime. Regarding utility, increasing the forgetting ratio led to a decrease in the AUC of the ML model retrained from scratch ($M_{\mathbf{D}_r}^O$) for both MLP and XGBoost. For SISA with MLP, there was almost no impact on utility, while utility decreased slightly for SISA with XGBoost. EUPG with MLP ($k = 10$ and $\epsilon = 0.5$) showed remarkable resilience, maintaining high utility scores even when the forgetting ratio reached 50%. With XGBoost and $k = 10$ the utility decreased slightly as the forgetting ratio increased, while with $\epsilon = 0.5$, the performance was already poor for low forgetting ratios and deteriorated further as the forgetting ratio increased.

The runtime tended to decrease with increasing forgetting ratios for retraining from scratch and SISA, due to the lower training data size that requires less computational effort. EUPG also showed large reductions in its already very short runtime with increasing forgetting ratios.

5.4 Detailed runtime

Table 9 reports the runtimes of the k -anonymization process to obtain \mathbf{D}^k , train M^k , and fine-tune M^k on \mathbf{D} with different values of k . Note that the anonymization time decreases almost linearly as k increases, demonstrating that larger values of k reduce the anonymization effort due to more uniform (aggregated) training data. On the other hand, the training and fine-tuning times for both MLP and XGBoost models remain relatively stable across k values. This indicates that the dominant factor in runtime is the initial k -anonymization process rather than the subsequent training or fine-tuning steps.

Tables 10 and 11 report runtimes for DP anonymization on tabular and image data, respectively. Runtimes are given for the embedding process (for categorical attributes), to create \mathbf{D}^ϵ , to train M^ϵ , and to fine-tune M^ϵ on \mathbf{D} with different values of ϵ . In our experiments, we only required embeddings for the Adult data set. As expected, the runtime to generate ϵ -differentially private data sets (\mathbf{D}^ϵ) remained stable across different values of ϵ , as did the runtimes to train and fine-tune all models. This stability suggests that

Table 9: Runtime in seconds when using k -anonymity with different values of k

data set	k	k -Anonymizing \mathbf{D}	Training M^k		Fine-tuning M^k on \mathbf{D}	
			MLP	XGBoost	MLP	XGBoost
Adult income	3	161.97	93.67	5.06	4.93	0.41
	5	101.97	94.56	5.73	4.99	0.44
	10	50.59	94.13	5.13	4.86	0.42
	20	25.83	94.90	4.51	4.80	0.39
	80	6.94	94.31	2.88	4.92	0.38
Heart disease	3	90.69	195.12	2.33	5.08	0.23
	5	57.50	193.46	2.40	5.06	0.24
	10	28.90	193.18	2.33	4.98	0.24
	20	14.57	193.24	2.24	5.08	0.25
	80	3.67	195.13	1.26	5.07	0.19
Credit	3	356.47	440.86	3.48	11.44	0.47
	5	162.00	442.65	4.07	11.41	0.48
	10	102.83	442.26	3.33	11.17	0.51
	20	40.56	442.50	4.00	11.35	0.50
	80	10.24	442.58	3.34	11.30	0.50

the computational cost of implementing differential privacy through adjustments in ϵ is almost invariant to the choice of ϵ . An important cost is the initial data embedding or transformation process needed for data sets with semantically rich categorical attributes.

6 Conclusions and future work

We have introduced EUPG, a novel framework for efficient machine unlearning that guarantees privacy to individuals who requested data forgetting while maintaining the utility (accuracy) of the ML model. EUPG adopts formal privacy models for fast and private unlearning, which has the advantage over approximate unlearning assessed through MIAs of being adaptable to protecting against future attack developments beyond MIAs.

Our experiments on a variety of data sets, involving numerical and categorical tabular data as well as images, confirm that our approach retains the utility of the original ML model with less computational and storage costs than exact unlearning methods and with privacy guarantees that comply with the right to be forgotten.

Future work will include on the one side optimizing the pre-training and fine-tuning processes to reduce their computational cost, and on the other side adapting EUPG to other ML tasks and more complex ML models —such as large language models (LLMs) and other foundation models. Dealing with LLMs raises the challenges of applying privacy models to unstructured text data and ensuring that the method is scalable w.r.t. the size of the training data. Whereas in this paper we enforced the privacy model on the training data (pre-protection), enforcing it on the ML model parameters during training (in-protection) might be preferable for LLMs and foundation models, due to the huge size of the training data involved. Nevertheless, the main ideas proposed here (forgetting by falling back on the protected ML model and then fine-tuning on the retain data) should still be valid.

Table 10: Runtime in seconds when using DP with different ϵ values on tabular data

data set	ϵ	Embedding	Generating D^ϵ	Training M^ϵ		Fine-tuning M^ϵ on D	
				MLP	XGBoost	MLP	XGBoost
Adult income	0.5	167.43	21.17	93.97	6.52	4.60	0.42
	2.5		20.92	93.88	6.65	4.79	0.43
	5		20.90	94.69	6.79	4.60	0.42
	25		21.05	94.56	6.42	4.67	0.43
	50		20.90	94.55	6.41	4.91	0.43
	100		19.88	94.51	6.51	5.11	0.43
Heart disease	0.5	0	28.92	193.39	2.16	5.08	0.25
	2.5		28.94	194.82	2.37	5.18	0.25
	5		29.15	193.69	2.23	5.19	0.26
	25		28.88	193.83	2.36	4.96	0.27
	50		29.11	194.05	2.35	5.10	0.25
	100		29.01	195.07	2.45	5.05	0.28
Credit	0.5	0	43.43	442.56	3.79	12.03	0.32
	2.5		43.21	441.41	3.37	11.80	0.44
	5		43.29	440.94	3.89	11.59	0.44
	25		43.91	442.30	4.16	11.58	0.49
	50		43.62	442.15	4.20	11.45	0.48
	100		43.49	442.03	3.91	11.35	0.46

Table 11: Runtime in seconds when using DP with different ϵ values on image data

data set	ϵ	Embedding	Generating D^ϵ	Training M^ϵ	Fine-tuning M^ϵ on D
				DenseNet12	DenseNet12
CIFAR-10	0.5	0	95.41	28979.72	1662.23
	2.5		96.63	29028.00	1662.07
	5		96.05	29021.83	1665.11
	25		96.50	28939.69	1658.16
	50		96.32	29066.90	1707.61
	100		96.02	30044.02	1678.08

Acknowledgments

Partial support has been received from the Government of Catalonia (ICREA Acadèmia Prizes to J. Domingo-Ferrer and to D. Sánchez, and grant 2021SGR-00115), MCIN/AEI/ 10.13039/501100011033 and “ERDF A way of making Europe” under grant PID2021-123637NB-I00 “CURLING”, and the EU’s NextGenerationEU/PRTR via INCIBE (project “HERMES” and INCIBE-URV cybersecurity chair).

References

- [1] Sercan Ö Arik and Tomas Pfister. Tabnet: Attentive interpretable tabular learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 6679–6687, 2021.

- [2] Borja Balle, Giovanni Cherubin, and Jamie Hayes. Reconstructing training data with informed adversaries. In *2022 IEEE Symposium on Security and Privacy (SP)*, pages 1138–1156. IEEE, 2022.
- [3] Thomas Baumhauer, Pascal Schöttle, and Matthias Zeppelzauer. Machine unlearning: Linear filtration for logit-based classifiers. *Machine Learning*, 111(9):3203–3226, 2022.
- [4] Alberto Blanco-Justicia, Najeeb Jebreel, Benet Manzanares-Salor, David Sánchez, Josep Domingo-Ferrer, Guillem Collell, and Kuan Eeik Tan. Digital forgetting in large language models: A survey of unlearning methods. *Artificial Intelligence Review*, 58(3):90, 2025.
- [5] Lucas Bourtole, 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.
- [6] Jonathan Brophy and Daniel Lowd. Machine unlearning for random forests. In *International Conference on Machine Learning*, pages 1092–1104. PMLR, 2021.
- [7] Health Canada. Guidance document on public release of clinical information, 2019.
- [8] Jingyi Cao, Bo Liu, Yunqian Wen, Yunhui Zhu, Rong Xie, Li Song, Lin li, and Yaoyao Yin. Hiding among your neighbors: face image privacy protection with differential private k-anonymity. In *Proceedings of the 2022 IEEE International Symposium on Broadband Multimedia Systems and Broadcasting*, pages 1–6, 2022.
- [9] 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.
- [10] Nicholas Carlini, Steve Chien, Milad Nasr, Shuang Song, Andreas Terzis, and Florian Tramèr. Membership inference attacks from first principles. In *2022 IEEE Symposium on Security and Privacy (SP)*, pages 1897–1914. IEEE, 2022.
- [11] Min Chen, Zhikun Zhang, Tianhao Wang, Michael Backes, Mathias Humbert, and Yang Zhang. Graph unlearning. In *Proceedings of the 2022 ACM SIGSAC conference on computer and communications security*, pages 499–513, 2022.
- [12] Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pages 785–794, 2016.
- [13] Eli Chien, Haoyu Peter Wang, Ziang Chen, and Pan Li. Certified machine unlearning via noisy stochastic gradient descent. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [14] Eli Chien, Haoyu Peter Wang, Ziang Chen, and Pan Li. Langevin unlearning: A new perspective of noisy gradient descent for machine unlearning. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [15] Rishav Chourasia and Neil Shah. Forget unlearning: Towards true data-deletion in machine learning. In *International conference on machine learning*, pages 6028–6073. PMLR, 2023.
- [16] Josep Domingo-Ferrer and Vicenç Torra. Ordinal, continuous and heterogeneous k-anonymity through microaggregation. *Data Mining and Knowledge Discovery*, 11:195–212, 2005.

- [17] Cynthia Dwork. Differential privacy. In *International colloquium on automata, languages, and programming*, pages 1–12. Springer, 2006.
- [18] Cynthia Dwork. A firm foundation for private data analysis. *Communications of the ACM*, 54(1):86–95, 2011.
- [19] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *Theory of Cryptography: Third Theory of Cryptography Conference, TCC 2006, New York, NY, USA, March 4-7, 2006. Proceedings 3*, pages 265–284. Springer, 2006.
- [20] Cynthia Dwork and Aaron Roth. *The Algorithmic Foundations of Differential Privacy*. Now Publishers, 2014.
- [21] European Parliament and Council of the European Union. Regulation (EU) 2016/679 of the European Parliament and of the Council of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (General Data Protection Regulation). <https://data.europa.eu/eli/reg/2016/679/oj>, 2016.
- [22] Liyue Fan. Differential privacy for image publication. In *Theory and Practice of Differential Privacy (TPDP) Workshop*, volume 1, page 6, 2019.
- [23] Jonas Geiping, Hartmut Bauermeister, Hannah Dröge, and Michael Moeller. Inverting gradients-how easy is it to break privacy in federated learning? *Advances in neural information processing systems*, 33:16937–16947, 2020.
- [24] Antonio Ginart, Melody Guan, Gregory Valiant, and James Y Zou. Making ai forget you: Data deletion in machine learning. *Advances in neural information processing systems*, 32, 2019.
- [25] Aditya Golatkar, Alessandro Achille, Avinash Ravichandran, Marzia Polito, and Stefano Soatto. Mixed-privacy forgetting in deep networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 792–801, 2021.
- [26] Xinlei He, Hongbin Liu, Neil Zhenqiang Gong, and Yang Zhang. Semi-leak: Membership inference attacks against semi-supervised learning. In *European Conference on Computer Vision*, pages 365–381. Springer, 2022.
- [27] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017.
- [28] Hanxun Huang, Xingjun Ma, Sarah Monazam Erfani, James Bailey, and Yisen Wang. Unlearnable examples: Making personal data unexploitable. In *International Conference on Learning Representations*, 2021.
- [29] Ninghui Li, Tiancheng Li, and Suresh Venkatasubramanian. t -closeness: Privacy beyond k -anonymity and l -diversity. In *Proceedings of the IEEE 23rd International Conference on Data Engineering (ICDE 2007)*, pages 106–115, 2007.
- [30] Ashwin Machanavajjhala, Daniel Kifer, Johannes Gehrke, and Muthuramakrishnan Venkatasubramanian. l -diversity: privacy beyond k -anonymity. *ACM Transactions on Knowledge Discovery from Data*, 1(1):3, 2007.

- [31] Frank McSherry and Kunal Talwar. Mechanism design via differential privacy. In *48th Annual IEEE Symposium on Foundations of Computer Science (FOCS'07)*, pages 94–103. IEEE, 2007.
- [32] Siqiao Mu and Diego Klabjan. Rewind-to-delete: Certified machine unlearning for nonconvex functions. *arXiv preprint arXiv:2409.09778*, 2024.
- [33] Vasileios Mygdalis, Anastasios Tefas, and Ioannis Pitas. Introducing k-anonymity principles to adversarial attacks for privacy protection in image classification problems. In *In 2021 IEEE 31st International Workshop on Machine Learning for Signal Processing (MLSP)*, pages 1–6, 2021.
- [34] Taichi Nakamura, Yuiko Sakuma, and Hiroaki Nishi. Face-image anonymization as an application of multidimensional data k-anonymizer. *International Journal of Networking and Computing*, 11(1):102–119, 2021.
- [35] Thanh Tam Nguyen, Thanh Trung Huynh, Phi Le Nguyen, Alan Wee-Chung Liew, Hongzhi Yin, and Quoc Viet Hung Nguyen. A survey of machine unlearning. *arXiv preprint arXiv:2209.02299*, 2022.
- [36] European Parliament. Artificial intelligence act: European parliament legislative resolution of 13 march 2024 on the proposal for a regulation of the european parliament and of the council on laying down harmonised rules on artificial intelligence (artificial intelligence act) and amending certain union legislative acts (com(2021)0206 – c9-0146/2021 – 2021/0106(cod)). https://www.europarl.europa.eu/doceo/document/TA-9-2024-0138_EN.pdf, 2024. (Accessed on 04/30/2024).
- [37] Ahmed Salem, Yang Zhang, Mathias Humbert, Mario Fritz, and Michael Backes. MI-leaks: Model and data independent membership inference attacks and defenses on machine learning models. In *Network and Distributed Systems Security Symposium 2019*. Internet Society, 2019.
- [38] Pierangela Samarati. Protecting respondents identities in microdata release. *IEEE transactions on Knowledge and Data Engineering*, 13(6):1010–1027, 2001.
- [39] David Sánchez, Josep Domingo-Ferrer, Sergio Martínez, and Jordi Soria-Comas. Utility-preserving differentially private data releases via individual ranking microaggregation. *Information Fusion*, 30:1–14, 2016.
- [40] Sebastian Schelter, Stefan Grafberger, and Ted Dunning. Hedgecut: Maintaining randomised trees for low-latency machine unlearning. In *Proceedings of the 2021 International Conference on Management of Data*, pages 1545–1557, 2021.
- [41] Ayush Sekhari, Jayadev Acharya, Gautam Kamath, and Ananda Theertha Suresh. Remember what you want to forget: Algorithms for machine unlearning. *Advances in Neural Information Processing Systems*, 34:18075–18086, 2021.
- [42] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models. In *2017 IEEE symposium on security and privacy (SP)*, pages 3–18. IEEE, 2017.
- [43] Jordi Soria-Comas and Josep Domingo-Ferrer. Probabilistic k-anonymity through microaggregation and data swapping. In *2012 IEEE International Conference on Fuzzy Systems*, pages 1–8. IEEE, 2012.
- [44] Jordi Soria-Comas, Josep Domingo-Ferrer, David Sánchez, and Sergio Martínez. Enhancing data utility in differential privacy via microaggregation-based k-anonymity. *The VLDB Journal*, 23(5):771–794, 2014.

- [45] State of California Department of Justice, Office of the Attorney General. California consumer privacy act (ccpa). <https://oag.ca.gov/privacy/ccpa>, January 2018. Accessed: 2025-05-28.
- [46] Vinith Suriyakumar and Ashia C Wilson. Algorithms that approximate data removal: New results and limitations. *Advances in Neural Information Processing Systems*, 35:18892–18903, 2022.
- [47] Ayush K Tarun, Vikram S Chundawat, Murari Mandal, and Mohan Kankanhalli. Fast yet effective machine unlearning. *IEEE Transactions on Neural Networks and Learning Systems*, 2023.
- [48] Enayat Ullah, Tung Mai, Anup Rao, Ryan A Rossi, and Raman Arora. Machine unlearning via algorithmic stability. In *Conference on Learning Theory*, pages 4126–4142. PMLR, 2021.
- [49] Zihan Wang, Jason Lee, and Qi Lei. Reconstructing training data from model gradient, provably. In *International Conference on Artificial Intelligence and Statistics*, pages 6595–6612. PMLR, 2023.
- [50] Alexander Warnecke, Lukas Pirsch, Christian Wressnegger, and Konrad Rieck. Machine unlearning of features and labels. In *Proceedings 2023 Network and Distributed System Security Symposium*. Internet Society, 2023.
- [51] Lauren Watson, Chuan Guo, Graham Cormode, and Alexandre Sablayrolles. On the importance of difficulty calibration in membership inference attacks. In *International Conference on Learning Representations*, 2021.
- [52] Xiaokui Xiao and Yufei Tao. Anatomy: Simple and effective privacy preservation. In *Proceedings of the 32nd international conference on Very large data bases*, pages 139–150, 2006.
- [53] Jie Xu, Zihan Wu, Cong Wang, and Xiaohua Jia. Machine unlearning: Solutions and challenges. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 2024.
- [54] Haonan Yan, Xiaoguang Li, Ziyao Guo, Hui Li, Fenghua Li, and Xiaodong Lin. Arcane: An efficient architecture for exact machine unlearning. In *IJCAI*, volume 6, page 19, 2022.
- [55] Samuel Yeom, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. Privacy risk in machine learning: Analyzing the connection to overfitting. In *2018 IEEE 31st computer security foundations symposium (CSF)*, pages 268–282. IEEE, 2018.
- [56] Da Yu, Huishuai Zhang, Wei Chen, Jian Yin, and Tie-Yan Liu. How does data augmentation affect privacy in machine learning? In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 10746–10753, 2021.