
Deep Unlearn: Benchmarking Machine Unlearning

Abstract

1 Machine unlearning (MU) aims to remove the influence of particular data points
2 from the learnable parameters of a trained machine learning model. This is a
3 key capability in light of data privacy requirements, trustworthiness, and safety
4 in deployed models. MU is particularly challenging for deep neural networks
5 (DNNs), such as convolutional nets or vision transformers, as such DNNs tend to
6 memorize a notable portion of their training dataset. Nevertheless, the community
7 lacks a rigorous and multifaceted study that looks into the success of MU methods
8 for DNNs. In this paper, we investigate 18 state-of-the-art MU methods across
9 various benchmark datasets and models, with each evaluation conducted over
10 10 different initializations, a comprehensive evaluation involving MU over 100K
11 models. We show that, with the proper hyperparameters, Masked Small Gradients
12 (MSG) and Convolution Transpose (CT), consistently perform better in terms of
13 model accuracy and run-time efficiency across different models, datasets, and
14 initializations, assessed by population-based membership inference attacks (MIA)
15 and per-sample unlearning likelihood ratio attacks (U-LiRA). Furthermore, our
16 benchmark highlights the fact that comparing a MU method only with commonly
17 used baselines, such as Gradient Ascent (GA) or Successive Random Relabeling
18 (SRL), is inadequate, and we need better baselines like Negative Gradient Plus
19 (NG+) with proper hyperparameter selection.

20 1 Introduction

21 Machine unlearning aims to remove the influence of a specified subset of training data points from
22 trained models [10]. This process is crucial for enhancing privacy preservation, model safety, and
23 overall model quality. MU helps ensure compliance with the right to be forgotten [23], removes
24 erroneous data points that negatively impact model performance [43], and eliminates biases introduced
25 by parts of the training data [14]. Deep neural networks present significant challenges for MU due
26 to their computationally intensive training requirements, highly non-convex loss landscapes, and
27 tendency to memorize substantial portions of their training data [22, 12, 21]. The key open challenges
28 in MU include: (1) A degradation in model accuracy often accompanies unlearning; (2) Some MU
29 methods require the model to be trained in specific ways, such as saving checkpoints, tracking
30 accumulated gradients, or training with differential privacy, limiting their applicability to already
31 deployed models; (3) Assurance of information removal is difficult as there are no reliable metrics to
32 measure it accurately; (4) There is no consensus on which methods are the most effective.

33 A branch of MU known as *exact unlearning* aims to guarantee that the specified *forget* data have been
34 completely removed from the model. The most reliable exact MU method is to *retrain* the model
35 from scratch while excluding the forget data. Another exact MU that offers data removal guarantees is
36 SISA (Sharded, Isolated, Sliced, Aggregated) [9]. However, exact MU is computationally prohibitive,
37 emphasizing the need for more efficient methods. The alternative branch is *approximate unlearning*

38 that aims to approximate data deletion and is often less precise, but more computationally efficient.
39 Approximate MU lacks theoretical guarantees, necessitating empirical evaluations to determine their
40 effectiveness, reliability, and computational efficiency across various datasets. Thus, the community
41 needs a proper benchmarking and evaluation of state-of-the-art approximate MU methods.

42 In this paper, we benchmark 18 state-of-the-art approximate MU methods across 5 datasets and
43 2 DNN architectures commonly used in computer vision: ResNet18 [32] and TinyViT [48].
44 The 18 methods consist of 3 classical baseline methods Fine-tuning (FT), Gradient Ascent (GA),
45 Successive Random Labels (SRL); 7 high-ranking methods in the NeurIPS’2023 Machine Un-
46 learning competition [45], we name them Forget-Contrast-Strengthen (FCS), Masked-Small-
47 Gradients (MSG), Confuse-Finetune-Weaken (CFW), Prune-Reinitialize-Match-Quantize (PRMQ),
48 Convolution-Transpose (CT), Knowledge-Distillation-Entropy (KDE) and Repeated-Noise-Injection
49 (RNI); and 8 recently published papers: Saliency Unlearning (SalUN) [20], Catastrophic Forgetting-K
50 (CF-k) [24], Exact Unlearning-K (EU-k) [24], SCalable Remembering and Unlearning unBound
51 (SCRUB) [37], Bad Teacher (BT) [16], Fisher Forgetting (FF) [25], Influence Unlearning (IU) [35, 33]
52 and Negative Gradient Plus (NG+) [37]. We evaluate the different unlearning methods in terms of four
53 major aspects: privacy evaluation, performance retention, computational efficiency, and unlearning
54 reliability.

55 The **contributions** of our study is to address the following research questions:

56 *Q1. Are the commonly used MU baselines reliable?* Most of the recently introduced MU methods are
57 compared only with three baselines: Fine-tuning, Gradient Ascent, and Successive Random Labels,
58 and not across recently introduced approaches. We show that with proper hyperparameter selection
59 FT is a reliable baseline. In contrast, GA consistently performs poorly and should be replaced by the
60 more recent Negative Gradient Plus (NG+) [37] that simultaneously reduces the performance in the
61 forget set while maintaining the performance on the rest of the data points.

62 *Q2. How reliable are MU methods across datasets, and, models, and initializations?* Our findings show
63 that, unlike the majority of recent MU methods, Masked-Small-Gradients is consistently among the
64 best performing methods across various metrics. In contrast, methods such as SRL do not consistently
65 outperform others across different datasets.

66 *Q3. Among existing methods, which are the most reliable and accurate?* Among 18 methods,
67 our evaluation shows that certain methods, such as Masked-Small-Gradients (MSG), Convolution-
68 Transpose (CT), and Fine-tuning (FT), exhibit desirable properties. Specifically, MSG and CT show
69 resilience against U-LIRA, a strong per-sample membership inference attack. Furthermore, all three
70 methods show consistency across datasets, and these methods could serve as reliable baselines for
71 future studies.

72 We publish the source code used for reproducing the experiments conducted in this paper at [released
73 upon publishing this paper].

74 2 Background of Machine Unlearning

75 **Setting.** Starting with a **training set** $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ and a trained model f_O , referred to as the
76 **original model**, the objective of a MU method U is to *remove* the influence of a particular subset
77 of training set $\mathcal{D}_F = \{(\mathbf{x}_i, y_i)\}_{i=1}^K \subset \mathcal{D}$ referred to as the **forget set** where $K \ll N$. The rest of
78 training set $\mathcal{D}_R = \mathcal{D} \setminus \mathcal{D}_F$ is called the **retain set**. The forget set and the retain set are distinct and
79 complementary subsets of the training set. The outcome of MU is an **unlearned** model f_U , the aim
80 for which is to *perform* on par with a model *retrained from scratch* on \mathcal{D}_R ; this latter model f_R is
81 referred to as the **retrained** model. We denote the weights of the original model and the retrained
82 model as θ_O and θ_R , respectively. For evaluations, we consider two held-out sets: the **validation set**
83 \mathcal{D}_V and **test set** \mathcal{D}_T , both drawn from the same distribution as \mathcal{D} . We consider the accuracy of the
84 retrain model as the optimal accuracy. One critical assumption we make is that the MU method has
85 access to $\theta_O, \mathcal{D}_R, \mathcal{D}_F$, and \mathcal{D}_V .

86 **(1) Unlearning Evaluation.** To evaluate the success of unlearning, one approach is to check whether
 87 data points in \mathcal{D}_F still influence the predictions made by the unlearned model [13, 37, 31]. This is
 88 commonly done via *influence functions* [35, 8, 29], membership inference attacks (MIA) [44] MIA
 89 has become one of the most common approaches for evaluating unlearning algorithms. It aims to
 90 determine whether specific data points were part of the original training dataset based on the unlearned
 91 model. We compute the population-based U-MIA, denoted with $\text{MIA}(\mathcal{D}, f_U)$ evaluated on data \mathcal{D} ¹.
 92 Using this we define the **discernibility** metric as $\text{Disc}(\mathcal{D}_V, f_U) = |2 \times \text{MIA}(\mathcal{D}, f_U) - 1| \in [0, 1]$,
 93 and similarly, **indiscernibility** is given by $\text{Indisc}(\mathcal{D}, f_U) = 1 - \text{Disc}(\mathcal{D}, f_U) \in [0, 1]$. We set \mathcal{D} to
 94 either the test set \mathcal{D}_T or validation set \mathcal{D}_V . The indiscernibility equals 1 when the accuracy of the
 95 MIA is not better than random guessing. A more recent MIA variation is the unlearning likelihood
 96 ratio attack (U-LiRA) [37]. As highlighted by [31], U-LiRA is a more robust evaluation approach for
 97 approximate MU. Nonetheless, U-LiRA is much more computationally demanding than U-MIA. We
 98 evaluate the methods that defeat the weaker, less expensive U-MIA attack, additionally against the
 99 U-LiRA attack.

100 **(2) Accuracy.** The classification accuracy of unlearned model on the retain set should be as close
 101 as possible to that of the original model. First, we consider three metrics derived directly from the
 102 model’s classification accuracy on different sets: **retain accuracy** (RA), **forget accuracy** (FA), and
 103 **test accuracy** (TA). RA is defined as

$$\text{RA}(\mathcal{D}_R, f_U) = \frac{1}{|\mathcal{D}_R|} \sum_{(\mathbf{x}_i, y_i) \in \mathcal{D}_R} \mathbf{1}_{y_i = f_U(\mathbf{x}_i)} \in [0, 1]. \quad (1)$$

104 The metrics for FA and TA can be derived by replacing \mathcal{D}_R with \mathcal{D}_F and \mathcal{D}_T , respectively. Second,
 105 while RA, FA, and TA give us insight into the overall accuracy of the unlearned model, they do not
 106 capture how well it performs compared to a retrained model f_R . Considering the f_R model as a
 107 gold standard, we derive three more metrics: **retain retention** (RR), **forget retention** (FR), **test**
 108 **retention** (TR), where RR is given by,

$$\text{RR}(f_U, f_R) = \frac{\text{RA}(\mathcal{D}_R, f_U)}{\text{RA}(\mathcal{D}_R, f_R)} \in [0, +\infty), \quad (2)$$

109 and formulas for FR and TR can be derived using FA and TA, respectively. An unlearned model with
 110 a score of 1 indicates that its accuracy perfectly matches the accuracy of the reference retrain model.
 111 A score below 1 indicates that the model under-performs and a score above 1 indicates that the model
 112 over-performs. We further define the **retention deviation** (RetDev) as:

$$\text{RetDev} = |\text{RR}(f_U, f_R) - 1| + |\text{FR}(f_U, f_R) - 1| + |\text{TR}(f_U, f_R) - 1| \in [0, +\infty), \quad (3)$$

113 which provides information on the cumulative divergence of the unlearned model in terms of retention
 114 score. The closer to 0, the better as 0 indicates that the model perfectly matches the performance of
 115 the retrained model.

116 **(3) Efficiency.** run-time efficiency (RTE) of a MU method should ideally be lower than the naive
 117 approach of just retraining the model from scratch on the retain set. As the high computational cost
 118 of the retraining algorithm motivated the development of approximate MU methods, we evaluate how
 119 much faster each MU method is compared to the retraining algorithm. We define RTE of U as the
 120 number of seconds it takes to complete, denoted as $\text{RT}(U)$ (we use the same machine and resources
 121 for all the experiments). To indicate the relative speedup compared to the retrained model, we define
 122 RTE of an unlearn method U as:

$$\frac{\text{RT}(U_R)}{\text{RT}(U)} \in [0, +\infty), \quad (4)$$

123 where U_R denotes the retrain method. The RTE of retraining from scratch would thus be 1; any
 124 method with an RTE less than 1 is slower than retraining from scratch, and vice versa.

¹The methodology for this follows the classic MIA: we compute the losses on \mathcal{D}_F and \mathcal{D}_V , we shuffle and trim them so that they are of equal size. We then train logistic regression models in a 10-fold cross validation, and compute the average accuracy across the folds.

125 **3 Machine Unlearning Methods**

126 We briefly discuss the main unlearning methods considered.

127 **3.1 Classical Baselines**

128 **FineTune (FT)** finetunes the original model f_O on only the retain set \mathcal{D}_R for several epochs.

129 **Successive Random Labels (SRL)** the model is trained on both the forget set \mathcal{D}_F and \mathcal{D}_R where the
130 labels of \mathcal{D}_F are randomly assigned at each epoch.

131 **Gradient Ascent (GA)** trains the model using gradient *ascent* steps on the \mathcal{D}_F .

132 **3.2 State-of-the-art MU methods**

133 Expanding upon the classical baselines, we additionally evaluate 15 recent MU methods. We first
134 discuss the seven top-performing methods from the Machine Unlearning Competition 2023 on
135 Kaggle.

136 **Forget-Contrast-Strengthen (FCS)** [1] minimizes the Kullback-Leibler Divergence (KLD) between
137 the model’s output on \mathcal{D}_F and a uniform distribution over the output classes, then alternatively
138 optimizes a contrastive loss between the model’s outputs on \mathcal{D}_R and \mathcal{D}_F , and minimizes the cross-
139 entropy loss on \mathcal{D}_R .

140 **Masked-Small-Gradients (MSG)**[2] accumulates gradients via gradient *descent* on the \mathcal{D}_R and
141 gradient *ascent* on the \mathcal{D}_F , then reinitialize weights with the smallest absolute gradients while
142 dampening subsequent weights updates on the \mathcal{D}_R for the other weights.

143 **Confuse-Finetune-Weaken (CFW)**[3] injects noise into the convolutional layers and then trains the
144 model using a class-weighted cross-entropy on \mathcal{D}_R , then injects noise again toward the final epochs.

145 **Prune-Reinitialize-Match-Quantize (PRMQ)** [4] first prunes the model via L1 pruning, reinitializes
146 parts of the model, optimises it using a combination of cross-entropy and a mean-squared-error on the
147 entropy between the outputs of f_O and f_U on \mathcal{D}_R and finally converts f_U ’s weights to half-precision
148 floats.

149 **Convolution-Transpose** [5] simply transposes the weights in the convolutional layers and trains on
150 \mathcal{D}_R .

151 **Knowledge-Distillation-Entropy (KDE)** [6] uses a teacher-student setup. Both student and teacher
152 start as copies of the original model, then the student’s first and last layers are re-initialised. The
153 student f_U minimizes its Kullback-Leibler Divergence (KLD) with the f_O over \mathcal{D}_V , then minimizes
154 a combination of losses: a soft cross-entropy loss between f_U and f_O , a cross-entropy loss on outputs
155 of \mathcal{D}_R from f_U , and the KLD between f_U and f_O on \mathcal{D}_R .

156 **Repeated-Noise-Injection (RNI)** [7] first reinitialises the final layer of the model, then repeatedly
157 injects noise in different layer of the model while training on the \mathcal{D}_R .

158 We further consider eight state-of-the-art methods introduced in the literature.

159 **Fisher Forgetting (FF)** [25, 20] adds noise to f_O with zero mean and covariance determined by the
160 4th root of Fisher Information matrix with respect to θ_O on \mathcal{D}_R .

161 **Influence Unlearning (IU)** [33, 47, 34] uses Influence Functions[18] to determine the change in θ_O
162 if a training point is removed from the training loss. IU estimates the change in model parameters
163 from θ_O to the model trained without a given data point. We use the first-order WoodFisher-based
164 approximation from [34].

165 **Catastrophic Forgetting - K (CF-K)** [24] freezes the first layers then trains the last k layers of the
166 model on \mathcal{D}_R .

167 **Exact Unlearning - K (EU-K)** [24] freezes the first layers then restores the weights of the last k
 168 layers to their initialization state. We randomly reinitialize the weights instead, so that the method no
 169 longer requires knowledge about the training process of f_O .

170 **SCRUB** [37] leverages a student-teacher setup where the model is optimised for three objectives:
 171 matching the teacher’s output distribution on \mathcal{D}_R , correctly predicting the \mathcal{D}_R set and ensuring the
 172 output distributions of the teacher and student diverge on the \mathcal{D}_F

173 **Saliency Unlearning (SaLUN)** [20] determines via gradient *ascent* which weights of θ_O are the
 174 most relevant to \mathcal{D}_F , then trains the model simultaneously on \mathcal{D}_R and \mathcal{D}_F with random labels on
 175 \mathcal{D}_F , while dampening the gradient propagation based on the selected weights.

176 **Negative Gradient Plus (NG+)** [37] is an extension of the Gradient Ascent approach where
 177 additionally a gradient descent step is taken over the \mathcal{D}_R .

178 **Bad Teacher (BT)** [16] uses a teacher-student approach with two teachers: the original model, and a
 179 randomly initialized model - the bad teacher-, the student starts as copy of f_U then learns to mimic
 180 the f_O on \mathcal{D}_R and the bad teacher on the \mathcal{D}_F .

181 4 Experimental Evaluation

182 **Experiments.** We evaluate the 18 recent MU methods as described in Section 3 across 5 benchmark
 183 datasets: MNIST [38], FashionMNIST [49], CIFAR-10 [36], CIFAR-100 [36], and UTK-Face [53].
 184 These datasets vary in difficulty, number of classes, instances per class, and image sizes. We consider
 185 two model architectures: a TinyViT and a ResNet18 model. Hence, in total we evaluated nine
 186 different combinations of models and architectures: ResNet18 and TinyViT on MNIST, Fashion-
 187 MNIST, CIFAR-10, CIFAR-100, and ResNet18 on UTKFace. More information on the data sets,
 188 hyperparameters, and data augmentations used to train the original and retrained models is provided
 189 in the appendix B. *We construct the forget set by sampling 10% of \mathcal{D} .*

190 The performance of the MU methods can change across datasets, model configurations, and model
 191 initializations; a reliable MU method remains consistent across these changes. For each method,
 192 dataset and model combination, we unlearn from Original models initialized using 10 different seeds
 193 and consider the average performance across seeds.

194 A further observation is that prior research tends to compare MU methods with default hyperparamete-
 195 rs, potentially leading to a less competitive performance of the method. To ensure that we get the
 196 best performance out of each method, we perform three hyperparameter sweeps to find the best set of
 197 hyperparameters for each method. To ensure a fair comparison, we use same number of searches for
 198 each method. Each hyper-parameter sweep uses 100 trials to minimize four loss functions: Retain
 199 Loss ($\mathcal{L}_{\text{Retain}}$), Forget Loss ($\mathcal{L}_{\text{Forget}}$), Val Loss (\mathcal{L}_{Val}), and Val MIA ($\mathcal{L}_{\text{Val-MIA}}$) given by

$$\mathcal{L}_{\text{Retain}} = \alpha \times |\text{RA}(f_U) - \text{RA}(f_R)|, \quad \mathcal{L}_{\text{Forget}} = \beta \times |\text{FA}(f_U) - \text{FA}(f_R)|, \quad (5)$$

$$\mathcal{L}_{\text{Val}} = \gamma \times |\text{VA}(f_U) - \text{VA}(f_R)|, \quad \mathcal{L}_{\text{Val-MIA}} = \eta \times \text{Disc}(\mathcal{D}_V, f_U), \quad (6)$$

200 where the $\mathcal{L}_{\text{Retain}}$ captures the divergence in accuracy between the retrained and unlearned model
 201 over the \mathcal{D}_R , $\mathcal{L}_{\text{Forget}}$ and \mathcal{L}_{Val} capture the divergence over $\mathcal{D}_F, \mathcal{D}_V$ respectively and $\mathcal{L}_{\text{Val-MIA}}$ cap-
 202 tures whether the loss distributions over \mathcal{D}_F and \mathcal{D}_V are distinguishable from one another via the
 203 discernibility score defined in Section 2. We set $\alpha = \beta = \gamma = \frac{1}{3}$ and $\eta = 1$ as we found these
 204 values to balance the importance of importance retention and the resilience to Membership Inference
 205 Attacks. Per unlearn method, we use the hyperparameter configuration that minimises the four loss
 206 terms when evaluating the method. Thus, for each unlearning method, we first unlearn 300 models to
 207 do the hyper-parameter sweep, then unlearn 10 models with the best set of hyper-parameters, leading
 208 to 5, 580 per dataset for a given architecture, leading to a total of 50, 220 for the 9 dataset / model
 209 combinations.

210 **Ranking.** A challenge in the comparison of MU method performance comes from the potential
 211 proximity of the evaluation metrics. As a simple example, suppose we have four methods U_1, \dots, U_4

Table 1: Ranking by performance on Retention Deviation and Indiscernibility across datasets and architectures. We count the number of times each method appears in the Best Performers group (G1), Average performance group (G2) and Worst performers group (G3) (see §4). The final rank is computed based on the number of times the method appears in G1—with occurrences in G2 and G3 used to break ties if needed. If a method does not produce any usable models, it is assigned to a Failed group (F). Three methods appear in the top 3 for both performance measures: MSG (1st and 1st), CT (3rd and 1st) and KDE (3rd and 2nd).

Rank	Method	Retention Deviation				Rank	Method	Indiscernibility			
		G1	G2	G3	F			G1	G2	G3	F
1	FT	8	1	0	0	1	CT	9	0	0	0
1	MSG	8	1	0	0	1	MSG	9	0	0	0
2	PRMQ	7	2	0	0	2	CFW	7	2	0	0
3	CT	7	1	1	0	2	RNI	7	2	0	0
3	KDE	7	1	1	0	2	KDE	7	2	0	0
3	CFW	7	1	1	0	3	FT	6	3	0	0
4	FCS	6	3	0	0	3	PRMQ	6	3	0	0
4	SalUN	6	3	0	0	3	SalUN	6	3	0	0
5	NG+	5	4	0	0	4	SRL	6	2	1	0
5	SRL	5	4	0	0	5	NG+	5	4	0	0
6	SCRUB	4	3	1	1	5	FCS	5	4	0	0
7	BT	2	7	0	0	6	SCRUB	5	3	0	1
7	RNI	2	7	0	0	7	BT	4	5	0	0
8	CF-k	2	3	2	2	8	CF-k	1	2	4	2
9	IU	1	0	2	6	9	EU-k	1	2	2	4
10	EU-k	0	5	0	4	10	GA	0	4	4	1
11	GA	0	1	7	1	11	IU	0	0	3	6
12	FF	0	0	0	9	12	FF	0	0	0	9

with accuracies: 98%, 99%, 50%, 1%, respectively; if we simply rank the methods, the rank itself would not be representative of the fact that e.g. U_1 and U_2 are much above U_3 and U_4 . In order to enable distinctions based on proximities, we use Agglomerative Clustering and define cut-off points such that we obtain three clusters: (1) Best performers (G1), (2) Average performers (G2), and (3) Worst performers (G3). If a method does not produce 10 usable models, one per original model, it is assigned to a Failed group (F). For each method, we count the number of times it appears in each of the three groups (with nine being the maximum). To obtain a final ranking of the methods, we first rank the methods using the number of times it appears in the Best Performers group (G1); if ties occur, we use the Average Performers (G2) group to break them. If ties persist, the Worst Performers (G3) group serves as the final tie-breaker. This method ensures a clear and fair ranking by considering each performance group in order of importance.

5 Main Results

Table 1 presents the main results of our evaluations on MU methods based on Retention Deviation and Indiscernibility. The results for the run-time efficiency are shown in Table 2.

On the reliability of baselines. The commonly used baseline FT trains the original model only on the retain set for several epochs to enable the model to forget information about the forget set. In our evaluation, FT performs best based on the Retention Deviation, and is ranked third based on Indiscernibility. The latter observation may come from the fact that FT does not explicitly unlearn the forget set or perturb the model parameters. Based on these results we conclude it is a reasonable baseline to evaluate against. We however remark that since the mechanism underlying FT (training on the Retain set to maintain performance) is common to many other methods, these methods may inherit its susceptibility to MIA. Another common baseline, GA which performs gradient ascent on

Table 2: Run Time Efficiency on ResNet for the top performing methods. CT is the fastest on average, MSG runs up to 5x faster than naive retraining.

Unlearner	CIFAR-10	CIFAR-100	MNIST	FashionMNIST	UTKFace	Average
MSG	6.80	4.49	4.29	3.32	7.57	5.29
CFW	4.67	6.17	4.29	4.90	5.54	5.11
PRMQ	4.93	4.34	3.77	3.77	5.88	4.54
CT	17.49	11.82	5.83	4.47	13.34	10.59
KDE	6.33	3.98	3.27	3.22	8.19	5.00
FT	8.15	5.16	5.07	4.29	5.78	5.69

234 the forget set, performs poorly across both metrics. Its more recent variation NG+, which uses an
 235 additional retain set correction, ranks fifth for both metrics, making it a more suitable baseline.

236 **On the reliability of newly proposed unlearning methods.** MSG obtains a first rank in both
 237 Performance Retention Deviation and Indiscernibility. Its unique approach identifies the parameters
 238 in the Convolutional layers that most contribute to the information to be forgotten. This strategy,
 239 differing from FT, allows MSG to retain performance from the Retain set while modifying the weights
 240 that are more relevant to the Forget set. PRMQ ranks second in performance retention, making it one
 241 of the top performers. However, it suffers from the same lower performance in Indiscernibility as FT.
 242 PRMQ however does not leverage the Forget set; instead, it performs a form of knowledge distillation
 243 by attempting to reproduce the results of the original models on the Retain set. Additionally, during
 244 the pruning phase, it reinitializes the weights for the MLP and Convolutional layers. CT ranks third
 245 in Performance Retention and first in Indiscernibility. It is interesting to note that CT and MSG are
 246 both consistently among the top performers in Indiscernibility, where FT performs poorly.

247 **On the robustness across architectures.** Critical for a good MU method is its ability to generalise
 248 across various DNN architectures. We conducted experiments with both ResNet18 and TinyViT.
 249 Methods such as CT, despite being tailored to Convolutional Neural Networks (CNN) models, still
 250 perform competitively when applied to Vision Transformers. Methods such as CT and MSG, while
 251 proposed for CNN layers, work well on Vision Transformer as one can leverage the 2D Convolutions
 252 used in Positional Encoders. We provide additional details on the ranking based on architectures in
 253 Appendix F.

254 **On the speed of the unlearn methods.** MSG, the best MU candidate runs 7.6x faster than
 255 retraining from scratch on UTKFace and 3.3x on FashionMNIST. On average MSG runs 5.3x faster
 256 than retraining. The fastest method is CT which achieves a speedup of 17.5x compared to Retraining.
 257 This speedup stems from its simple approach: transpose the weights of the convolutional layers.

258 **On the performance when evaluated against a stronger MIA.** From the results above, we
 259 note that there exist MU methods that perform fast and reliably across datasets, architectures, and
 260 initializations. Specifically, CT, FT, MSG seem strong candidates for reliable MU methods. However,
 261 as has also been highlighted in prior work [46], MIA has been debated as a strong metric, as its
 262 ability to assess MU is hampered by its own ability to infer data membership. In line with this,
 263 recent works have introduced more powerful variations on MIA [51] with [37] proposing the stronger
 264 U-LiRA attack for MU. For the best performers in Table 1, we apply the attack setup from [31] where
 265 we generate a total of 640 models (with varying train, retain and forget sets) and for each unlearn
 266 method perform a hyperparameter sweep to find the best configuration (for details see Appendix
 267 D). We determine for each data point its U-LiRA Inference Accuracy and report for each method its
 268 average and standard deviation (see Figure 1). From this, we conclude that MSG as well as CT resist
 269 U-LiRA attacks. Both MSG and CT thus not only rank first in terms of Performance Retention and
 270 Indiscernibility based on U-MIA, but also are robust against a stronger variation of MIA.

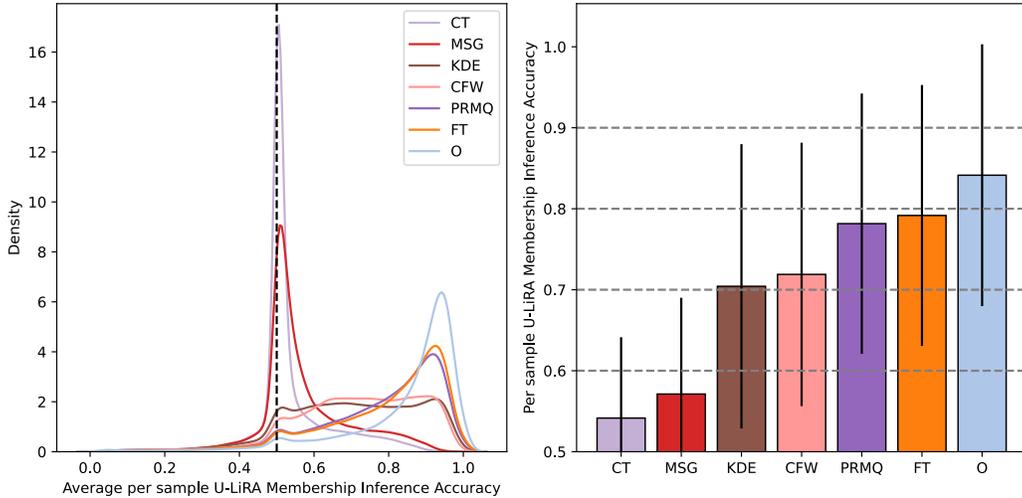


Figure 1: U-LiRA on CIFAR-10 on ResNet models. Both CT and MSG, which ranked first against U-MIA, showed great resilience against the U-LiRA attack.

271 6 Discussion and Conclusion

272 The increasing focus on data privacy and trustworthiness of machine learning models underscores the
 273 need for robust and practical methods to unlearn and remove the influence of specific data from trained
 274 models. Due to the growing size of models, we require methods that avoid the computationally costly
 275 retraining from scratch. In this work we performed a comprehensive comparison of approximate
 276 unlearning methods across various models and datasets aimed to address this critical issue.

277 We experimentally compared 18 methods across different datasets and architectures, focusing on
 278 assessing the method’s ability to maintain privacy and accuracy while being computationally efficient
 279 and reliable across datasets, architectures and random seeds. Our findings indicate that Masked-Small-
 280 Gradients, which accumulates gradients via gradient descent on the data to remember and gradient
 281 ascent on the data to forget to determine which weights to update, consistently outperforms for all
 282 metrics across the studied datasets, architectures, and initialization seeds. Similarly, Convolution
 283 Transpose, which leverages the simple transposition in convolutional layers, performed strongly.

284 Both CT and MSG were resistant against both a population-based Membership Inference Attack
 285 (MIA) and a stronger, per-sample attack (U-LiRA). However, a core challenge of approximate
 286 unlearning is that these methods will only be as strong as the attacks against which they are tested. As
 287 stronger and more complex attacks emerge, some approximate unlearning methods might no longer
 288 be as efficient as initially expected. This highlights the need for continuous evaluation and adaptation
 289 of unlearning methods to maintain their effectiveness. We also conducted experiments based on
 290 L2 distances, but found that no method consistently got close to the reference models’ weights, we
 291 provide further information in Appendix G.

292 **Limitations.** Due to computational costs, we limited our analysis to Tiny Vision Transformers and
 293 ResNet; a further investigation of other architectures could provide useful insights. We did not
 294 investigate different amounts of unlearning samples, which some methods are known to be sensitive
 295 to [37]. We did not consider repeated deletion, instead we assume that there is a single forget set and
 296 that the unlearning process happens once, as is common in the literature, nonetheless, in practical
 297 applications one might need to unlearn different smaller forget sets over time and some unlearning
 298 methods might not work as well under such scenario. We finally remark once again on the difficulty
 299 of evaluation for approximate unlearning [31]: while these methods provide significant gains in
 300 efficiency, novel attacks might highlight yet unknown weaknesses of the unlearning processes.

301 **Future work.** First, we put our focus on natural image data, however, machine unlearning is relevant
302 to other data types such as medical images or other modalities such as time series, audio and speech,
303 or language data. Second, we focus on the classification task, however, other learning tasks would
304 greatly benefit from machine unlearning too. For instance removing concepts from generative models
305 for images [20] or poisoned data in language models [31]. Third, this work focuses on empirically
306 benchmarking approximate machine unlearning methods. We do not provide a theoretical analysis of
307 these methods or a rigorous comparison with exact unlearning algorithms.

308 **Impact statement.** This paper aims to highlight the importance of effectively assessing approximate
309 machine unlearning methods. Our goal is to stress the need for evaluating new unlearning methods
310 against more reliable baselines and experimental setups. Additionally, it is crucial to assess the
311 consistency of a new unlearning method across various datasets and model architectures. Without
312 such a thorough evaluations, proposed unlearning methods may provide a false sense of privacy and
313 safety, ultimately limiting their effectiveness for data regulation.

314 References

- 315 [1] 2023. URL [https://www.kaggle.com/competitions/
316 neurips-2023-machine-unlearning/discussion/458721](https://www.kaggle.com/competitions/neurips-2023-machine-unlearning/discussion/458721). Accessed on January
317 31, 2024.
- 318 [2] 2023. URL [https://www.kaggle.com/competitions/
319 neurips-2023-machine-unlearning/discussion/459200](https://www.kaggle.com/competitions/neurips-2023-machine-unlearning/discussion/459200). Accessed on January
320 31, 2024.
- 321 [3] 2023. URL [https://www.kaggle.com/competitions/
322 neurips-2023-machine-unlearning/discussion/459334](https://www.kaggle.com/competitions/neurips-2023-machine-unlearning/discussion/459334). Accessed on January
323 31, 2024.
- 324 [4] 2023. URL [https://www.kaggle.com/competitions/
325 neurips-2023-machine-unlearning/discussion/459148](https://www.kaggle.com/competitions/neurips-2023-machine-unlearning/discussion/459148). Accessed on January
326 31, 2024.
- 327 [5] 2023. URL [https://www.kaggle.com/competitions/
328 neurips-2023-machine-unlearning/discussion/458531](https://www.kaggle.com/competitions/neurips-2023-machine-unlearning/discussion/458531). Accessed on January
329 31, 2024.
- 330 [6] 2023. URL [https://www.kaggle.com/competitions/
331 neurips-2023-machine-unlearning/discussion/458740](https://www.kaggle.com/competitions/neurips-2023-machine-unlearning/discussion/458740). Accessed on January
332 31, 2024.
- 333 [7] 2023. URL [https://www.kaggle.com/competitions/
334 neurips-2023-machine-unlearning/discussion/459095](https://www.kaggle.com/competitions/neurips-2023-machine-unlearning/discussion/459095). Accessed on January
335 31, 2024.
- 336 [8] Samyadeep Basu, Philip Pope, and Soheil Feizi. Influence Functions in Deep Learning Are
337 Fragile, February 2021. URL <http://arxiv.org/abs/2006.14651>. arXiv:2006.14651 [cs,
338 stat].
- 339 [9] Lucas Bourtole, Varun Chandrasekaran, Christopher A. Choquette-Choo, Hengrui Jia, Adelin
340 Travers, Baiwu Zhang, David Lie, and Nicolas Papernot. Machine Unlearning. In *2021 IEEE
341 Symposium on Security and Privacy (SP)*, pages 141–159, May 2021. doi: 10.1109/SP40001.
342 2021.00019. ISSN: 2375-1207.
- 343 [10] Yinzhi Cao and Junfeng Yang. Towards making systems forget with machine unlearning. In
344 *2015 IEEE symposium on security and privacy*, pages 463–480. IEEE, 2015.

- 345 [11] Yinzhi Cao and Junfeng Yang. Towards Making Systems Forget with Machine Unlearning.
346 In *2015 IEEE Symposium on Security and Privacy*, pages 463–480, San Jose, CA, May 2015.
347 IEEE. ISBN 978-1-4673-6949-7. doi: 10.1109/SP.2015.35. URL [https://ieeexplore.
348 ieee.org/document/7163042/](https://ieeexplore.ieee.org/document/7163042/).
- 349 [12] Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej Kos, and Dawn Song. The Secret
350 Sharer: Evaluating and Testing Unintended Memorization in Neural Networks, July 2019. URL
351 <http://arxiv.org/abs/1802.08232>. arXiv:1802.08232 [cs].
- 352 [13] Nicholas Carlini, Steve Chien, Milad Nasr, Shuang Song, Andreas Terzis, and Florian Tramèr.
353 Membership Inference Attacks From First Principles. In *2022 IEEE Symposium on Security
354 and Privacy (SP)*, pages 1897–1914, San Francisco, CA, USA, May 2022. IEEE. ISBN 978-
355 1-66541-316-9. doi: 10.1109/SP46214.2022.9833649. URL [https://ieeexplore.ieee.
356 org/document/9833649/](https://ieeexplore.ieee.org/document/9833649/).
- 357 [14] Ruizhe Chen, Jianfei Yang, Huimin Xiong, Jianhong Bai, Tianxiang Hu, Jin Hao, Yang
358 Feng, Joey Tianyi Zhou, Jian Wu, and Zuozhu Liu. Fast Model DeBias with Machine Un-
359 learning. *Advances in Neural Information Processing Systems*, 36:14516–14539, Decem-
360 ber 2023. URL [https://proceedings.neurips.cc/paper_files/paper/2023/hash/
361 2ecc80084c96cc25b11b0ab995c25f47-Abstract-Conference.html](https://proceedings.neurips.cc/paper_files/paper/2023/hash/2ecc80084c96cc25b11b0ab995c25f47-Abstract-Conference.html).
- 362 [15] Anna Choromanska, Mikael Henaff, Michael Mathieu, Gérard Ben Arous, and Yann LeCun.
363 The Loss Surfaces of Multilayer Networks, January 2015. URL [http://arxiv.org/abs/
364 1412.0233](http://arxiv.org/abs/1412.0233). arXiv:1412.0233 [cs].
- 365 [16] Vikram S. Chundawat, Ayush K. Tarun, Murari Mandal, and Mohan Kankanhalli. Can Bad
366 Teaching Induce Forgetting? Unlearning in Deep Networks Using an Incompetent Teacher.
367 *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(6):7210–7217, June 2023.
368 ISSN 2374-3468. doi: 10.1609/aaai.v37i6.25879. URL [https://ojs.aaai.org/index.
369 php/AAAI/article/view/25879](https://ojs.aaai.org/index.php/AAAI/article/view/25879). Number: 6.
- 370 [17] Vikram S. Chundawat, Ayush K. Tarun, Murari Mandal, and Mohan Kankanhalli. Zero-
371 Shot Machine Unlearning. *IEEE Transactions on Information Forensics and Security*, 18:
372 2345–2354, 2023. ISSN 1556-6013, 1556-6021. doi: 10.1109/TIFS.2023.3265506. URL
373 <https://ieeexplore.ieee.org/document/10097553/>.
- 374 [18] R. Dennis Cook and Sanford Weisberg. *Residuals and Influence in Regression*. New York:
375 Chapman and Hall, 1982. URL <http://pur1.umn.edu/37076>.
- 376 [19] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
377 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,
378 Jakob Uszkoreit, and Neil Houlsby. An Image is Worth 16x16 Words: Transformers for Image
379 Recognition at Scale. *arXiv:2010.11929 [cs]*, June 2021. URL [http://arxiv.org/abs/
380 2010.11929](http://arxiv.org/abs/2010.11929). arXiv: 2010.11929.
- 381 [20] Chongyu Fan, Jiancheng Liu, Yihua Zhang, Eric Wong, Dennis Wei, and Sijia Liu. SalUn:
382 Empowering Machine Unlearning via Gradient-based Weight Saliency in Both Image
383 Classification and Generation, March 2024. URL <http://arxiv.org/abs/2310.12508>.
384 arXiv:2310.12508 [cs].
- 385 [21] Vitaly Feldman and Chiyuan Zhang. What Neural Networks Memorize and Why: Discovering
386 the Long Tail via Influence Estimation, August 2020. URL [http://arxiv.org/abs/2008.
387 03703](http://arxiv.org/abs/2008.03703). arXiv:2008.03703 [cs, stat].
- 388 [22] Matt Fredrikson, Somesh Jha, and Thomas Ristenpart. Model Inversion Attacks that Exploit
389 Confidence Information and Basic Countermeasures. In *Proceedings of the 22nd ACM SIGSAC
390 Conference on Computer and Communications Security*, pages 1322–1333, Denver Colorado
391 USA, October 2015. ACM. ISBN 978-1-4503-3832-5. doi: 10.1145/2810103.2813677. URL
392 <https://dl.acm.org/doi/10.1145/2810103.2813677>.

- 393 [23] Antonio Ginart, Melody Guan, Gregory Valiant, and James Y Zou. Making ai forget you: Data
394 deletion in machine learning. *Advances in neural information processing systems*, 32, 2019.
- 395 [24] Shashwat Goel, Ameya Prabhu, Amartya Sanyal, Ser-Nam Lim, Philip Torr, and Ponnurangam
396 Kumaraguru. Towards Adversarial Evaluations for Inexact Machine Unlearning, February 2023.
397 URL <http://arxiv.org/abs/2201.06640>. arXiv:2201.06640 [cs].
- 398 [25] Aditya Golatkar, Alessandro Achille, and Stefano Soatto. Eternal Sunshine of the Spotless Net:
399 Selective Forgetting in Deep Networks, March 2020. URL <http://arxiv.org/abs/1911.04933>.
400 arXiv:1911.04933 [cs, stat].
- 401 [26] Aditya Golatkar, Alessandro Achille, and Stefano Soatto. Forgetting Outside the Box: Scrubbing
402 Deep Networks of Information Accessible from Input-Output Observations, October 2020. URL
403 <http://arxiv.org/abs/2003.02960>. arXiv:2003.02960 [cs, math, stat].
- 404 [27] Aditya Golatkar, Alessandro Achille, Avinash Ravichandran, Marzia Polito, and Stefano Soatto.
405 Mixed-Privacy Forgetting in Deep Networks, June 2021. URL [http://arxiv.org/abs/](http://arxiv.org/abs/2012.13431)
406 [2012.13431](http://arxiv.org/abs/2012.13431). arXiv:2012.13431 [cs].
- 407 [28] Laura Graves, Vineel Nagisetty, and Vijay Ganesh. Amnesiac Machine Learning, October 2020.
408 URL <http://arxiv.org/abs/2010.10981>. arXiv:2010.10981 [cs].
- 409 [29] Roger Grosse, Juhan Bae, Cem Anil, Nelson Elhage, Alex Tamkin, Amirhossein Tajdini,
410 Benoit Steiner, Dustin Li, Esin Durmus, Ethan Perez, Evan Hubinger, Kamilè Lukošiušė,
411 Karina Nguyen, Nicholas Joseph, Sam McCandlish, Jared Kaplan, and Samuel R. Bowman.
412 Studying Large Language Model Generalization with Influence Functions, August 2023. URL
413 <http://arxiv.org/abs/2308.03296>. arXiv:2308.03296 [cs, stat].
- 414 [30] Chuan Guo, Tom Goldstein, Awni Hannun, and Laurens van der Maaten. Certified Data Removal
415 from Machine Learning Models, August 2020. URL <http://arxiv.org/abs/1911.03030>.
416 arXiv:1911.03030 [cs, stat].
- 417 [31] Jamie Hayes, Iliia Shumailov, Eleni Triantafillou, Amr Khalifa, and Nicolas Papernot. Inexact
418 Unlearning Needs More Careful Evaluations to Avoid a False Sense of Privacy, March 2024.
419 URL <http://arxiv.org/abs/2403.01218>. arXiv:2403.01218 [cs].
- 420 [32] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image
421 Recognition. *arXiv:1512.03385 [cs]*, December 2015. URL [http://arxiv.org/abs/1512.](http://arxiv.org/abs/1512.03385)
422 [03385](http://arxiv.org/abs/1512.03385). arXiv: 1512.03385.
- 423 [33] Zachary Izzo, Mary Anne Smart, Kamalika Chaudhuri, and James Zou. Approximate Data
424 Deletion from Machine Learning Models, February 2021. URL [http://arxiv.org/abs/](http://arxiv.org/abs/2002.10077)
425 [2002.10077](http://arxiv.org/abs/2002.10077). arXiv:2002.10077 [cs, stat].
- 426 [34] Jinghan Jia, Jiancheng Liu, Parikshit Ram, Yuguang Yao, Gaowen Liu, Yang Liu, Pranay
427 Sharma, and Sijia Liu. Model Sparsity Can Simplify Machine Unlearning, January 2024. URL
428 <http://arxiv.org/abs/2304.04934>. arXiv:2304.04934 [cs].
- 429 [35] Pang Wei Koh and Percy Liang. Understanding Black-box Predictions via Influence Func-
430 tions, March 2017. URL <http://arxiv.org/abs/1703.04730>. arXiv:1703.04730 [cs, stat]
431 version: 1.
- 432 [36] Alex Krizhevsky. Learning Multiple Layers of Features from Tiny Images. *University of*
433 *Toronto*, 2009.
- 434 [37] Meghdad Kurmanji, Peter Triantafillou, Jamie Hayes, and Eleni Triantafillou. Towards Un-
435 bounded Machine Unlearning, October 2023. URL <http://arxiv.org/abs/2302.09880>.
436 arXiv:2302.09880 [cs].

- 437 [38] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document
438 recognition. *Proceedings of the IEEE*, 86(11):2278–2324, November 1998. ISSN 1558-2256.
439 doi: 10.1109/5.726791. Conference Name: Proceedings of the IEEE.
- 440 [39] Thanh Tam Nguyen, Thanh Trung Huynh, Phi Le Nguyen, Alan Wee-Chung Liew, Hongzhi
441 Yin, and Quoc Viet Hung Nguyen. A Survey of Machine Unlearning, October 2022. URL
442 <http://arxiv.org/abs/2209.02299>. arXiv:2209.02299 [cs].
- 443 [40] Sebastian Schelter. “Amnesia” – Towards Machine Learning Models That Can Forget User
444 Data Very Fast. In *Conference on Innovative Data Systems Research*, 2020.
- 445 [41] Sebastian Schelter, Stefan Grafberger, and Ted Dunning. HedgeCut: Maintaining randomised
446 trees for low-latency machine unlearning. In *Proceedings of the 2021 international conference
447 on management of data, SIGMOD ’21*, pages 1545–1557, New York, NY, USA, 2021. Associ-
448 ation for Computing Machinery. ISBN 978-1-4503-8343-1. doi: 10.1145/3448016.3457239.
449 URL <https://doi.org/10.1145/3448016.3457239>. Number of pages: 13 Place: Virtual
450 Event, China.
- 451 [42] Ayush Sekhari, Jayadev Acharya, Gautam Kamath, and Ananda Theertha Suresh. Remember
452 What You Want to Forget: Algorithms for Machine Unlearning, July 2021. URL <http://arxiv.org/abs/2103.03279>.
453 arXiv:2103.03279 [cs].
- 454 [43] Thanveer Shaik, Xiaohui Tao, Haoran Xie, Lin Li, Xiaofeng Zhu, and Qing Li. Exploring the
455 Landscape of Machine Unlearning: A Comprehensive Survey and Taxonomy, May 2023. URL
456 <http://arxiv.org/abs/2305.06360>. arXiv:2305.06360 [cs].
- 457 [44] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership Inference
458 Attacks Against Machine Learning Models. In *2017 IEEE Symposium on Security and Privacy
459 (SP)*, pages 3–18, San Jose, CA, USA, May 2017. IEEE. ISBN 978-1-5090-5533-3. doi:
460 10.1109/SP.2017.41. URL <http://ieeexplore.ieee.org/document/7958568/>.
- 461 [45] Eleni Triantafillou, Fabian Pedregosa, Jamie Hayes, Peter Kairouz, Isabelle Guyon, Meghdad
462 Kurmanji, Gintare Karolina Dziugaite, Peter Triantafillou, Kairan Zhao, Lisheng Sun Hosoya,
463 Julio C. S. Jacques Junior, Vincent Dumoulin, Ioannis Mitliagkas, Sergio Escalera, Jun Wan,
464 Sohier Dane, Maggie Demkin, and Walter Reade. Neurips 2023 - machine unlearning, 2023.
465 URL <https://kaggle.com/competitions/neurips-2023-machine-unlearning>.
- 466 [46] Yiwen Tu, Pingbang Hu, and Jiaqi Ma. Towards Reliable Empirical Machine Unlearning
467 Evaluation: A Game-Theoretic View, April 2024. URL [http://arxiv.org/abs/2404.
468 11577](http://arxiv.org/abs/2404.11577). arXiv:2404.11577 [cs].
- 469 [47] Alexander Warnecke, Lukas Pirch, Christian Wressnegger, and Konrad Rieck. Machine Un-
470 learning of Features and Labels, August 2023. URL <http://arxiv.org/abs/2108.11577>.
471 arXiv:2108.11577 [cs].
- 472 [48] Kan Wu, Jinnian Zhang, Houwen Peng, Mengchen Liu, Bin Xiao, Jianlong Fu, and Lu Yuan.
473 TinyViT: Fast Pretraining Distillation for Small Vision Transformers. In Shai Avidan, Gabriel
474 Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner, editors, *Computer
475 Vision – ECCV 2022*, volume 13681, pages 68–85. Springer Nature Switzerland, Cham, 2022.
476 ISBN 978-3-031-19802-1 978-3-031-19803-8. doi: 10.1007/978-3-031-19803-8_5. URL
477 https://link.springer.com/10.1007/978-3-031-19803-8_5. Series Title: Lecture
478 Notes in Computer Science.
- 479 [49] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MNIST: a Novel Image Dataset for
480 Benchmarking Machine Learning Algorithms, September 2017. URL [http://arxiv.org/
481 abs/1708.07747](http://arxiv.org/abs/1708.07747). arXiv:1708.07747 [cs, stat] version: 2.
- 482 [50] Heng Xu, Tianqing Zhu, Lefeng Zhang, Wanlei Zhou, and Philip S. Yu. Machine Unlearning:
483 A Survey, June 2023. URL <http://arxiv.org/abs/2306.03558>. arXiv:2306.03558 [cs].

- 484 [51] Jiayuan Ye, Aadyaa Maddi, Sasi Kumar Murakonda, Vincent Bindschaedler, and Reza Shokri.
485 Enhanced Membership Inference Attacks against Machine Learning Models. In *Proceedings*
486 *of the 2022 ACM SIGSAC Conference on Computer and Communications Security*, pages
487 3093–3106, Los Angeles CA USA, November 2022. ACM. ISBN 978-1-4503-9450-5. doi:
488 10.1145/3548606.3560675. URL <https://dl.acm.org/doi/10.1145/3548606.3560675>.
- 489 [52] Haibo Zhang, Toru Nakamura, Takamasa Isohara, and Kouichi Sakurai. A Review on Machine
490 Unlearning. *SN Computer Science*, 4(4):337, April 2023. ISSN 2661-8907. doi: 10.1007/
491 s42979-023-01767-4. URL <https://doi.org/10.1007/s42979-023-01767-4>.
- 492 [53] Zhifei Zhang, Yang Song, and Hairong Qi. Age Progression/Regression by Condi-
493 tional Adversarial Autoencoder, March 2017. URL <http://arxiv.org/abs/1702.08423>.
494 arXiv:1702.08423 [cs].

495 **Checklist**

- 496 1. For all authors...
- 497 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
- 498 contributions and scope? [Yes]
- 499 (b) Did you describe the limitations of your work? [Yes]
- 500 (c) Did you discuss any potential negative societal impacts of your work? [Yes]
- 501 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
- 502 them? [Yes]
- 503 2. If you are including theoretical results...
- 504 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 505 (b) Did you include complete proofs of all theoretical results? [N/A]
- 506 3. If you ran experiments...
- 507 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
- 508 mental results (either in the supplemental material or as a URL)? [Yes]
- 509 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
- 510 were chosen)? [Yes]
- 511 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
- 512 ments multiple times)? [Yes]
- 513 (d) Did you include the total amount of compute and the type of resources used (e.g., type
- 514 of GPUs, internal cluster, or cloud provider)? [Yes]
- 515 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 516 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 517 (b) Did you mention the license of the assets? [Yes]
- 518 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 519
- 520 (d) Did you discuss whether and how consent was obtained from people whose data you're
- 521 using/curating? [N/A]
- 522 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 523 information or offensive content? [N/A]
- 524 5. If you used crowdsourcing or conducted research with human subjects...
- 525 (a) Did you include the full text of instructions given to participants and screenshots, if
- 526 applicable? [N/A]
- 527 (b) Did you describe any potential participant risks, with links to Institutional Review
- 528 Board (IRB) approvals, if applicable? [N/A]
- 529 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 530 spent on participant compensation? [N/A]