# A MULTI-DOMAIN BENCHMARK FOR MACHINE UNLEARNING IN CLASSIFICATION TASKS

**Anonymous authors** 

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

025

026027028

029

031

033

034

035

037

040

041

042

043

044

045

046

047

048

051

052

Paper under double-blind review

#### **ABSTRACT**

Machine unlearning (MU), the process of removing specific data influences from trained machine learning models, is critical for regulatory compliance (e.g., GDPR's right to be forgotten) and for addressing copyright and privacy concerns in large-scale models. While a wide range of methods and metrics have been proposed, systematic evaluations remain fragmented, typically limited in scope by modality, metric coverage, or the number of methods considered. In this work, we present the most comprehensive MU benchmark to date, evaluating 12 unlearning methods on 8 datasets and models across four modalities (images, text, tabular data, and graphs) by assessing the three key aspects of an unlearning outcome: utility – the overall performance of the model after unlearning – efficacy – how well the data is forgotten – and *efficiency* – the computational cost of unlearning. We also introduce LUMA (Laplacian Unlearning Multidimensional Assessment), a unified metric that consolidates them into a single score. Unlike prior metrics, LUMA can flexibly incorporate multiple measures within each dimension (e.g., F1 over test and forget set for utility, UMIA for efficacy, runtime and GPU memory for efficiency), enabling more accurate and extensible comparisons. Our code is reproducible and extensible to serve as a benchmark for MU research.

#### 1 Introduction

The growing inclusion of machine learning (ML) models across various industries has raised concerns about using potentially sensitive data in model training (Grynbaum & Mac, 2023). In response, regulations such as the AI Act and the General Data Protection Regulation (GDPR) established the *right to be forgotten*, mandating that an individual's data must be removed upon request (Mantelero, 2013). In such cases, the model's owner must release a new version of the model itself, specifically excluding those targeted samples from the training set. However, retraining ML models from scratch upon every request is often too expensive in terms of time, money, and environmental costs (Crawford, 2022). *Machine unlearning* (MU) offers a promising alternative: instead of (re)training the model from scratch, MU aims to efficiently remove the influence of specific data points from an already trained model (Xu et al., 2024; Le Quy et al., 2022).

Despite its recent emergence, the literature on MU has been growing at a massive rate, outpacing benchmarking and surveying efforts. Specifically, while methods are often presented as generally applicable (Chundawat et al., 2023; Foster et al., 2024), most benchmarks focus on a single domain, such as images (see Section 2), which limits the ability to assess their generalizability. The Tabular and Graph domains remain largely unexplored: to the best of our knowledge, ours is the first MU benchmark covering these two domains. In contrast, the textual domain has received more attention; however, existing benchmarks target the text generation task (Maini et al., 2024; Shi et al., 2024) rather than the classification one.

In this work, we establish coherence in this fragmented landscape by studying MU methods on diverse domains under a single, unified evaluation benchmark. As illustrated in Figure 1, the proposed unlearning benchmark consists of three steps: (1) selecting datasets from one of the four supported domains, (2) training the appropriate model (between two sizes), and (3) applying unlearning methods directly to the model without requiring domain-specific adjustments.

After unlearning, step (4) is the collection and analysis of the values of the evaluation metrics. As detailed in Hayes et al. (2024), MU consists of three key aspects that must be evaluated jointly:

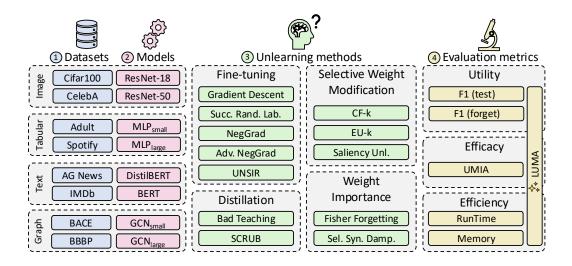


Figure 1: Experimental workflow of our benchmark. We evaluate 8 classification datasets across 4 domains, training 2 models per domain. We test 12 unlearning methods and assess them across 3 evaluation dimensions. We introduce a unified metric, LUMA, to facilitate comparison.

utility – how well the model performs after unlearning – efficacy – the degree to which the target data is removed – and efficiency – the computational cost of unlearning w.r.t. retraining the model from scratch. After presenting each of the three aspects separately, we introduce LUMA, a unified metric designed to bridge the current gap in the literature.

Our contributions are the following: (i). We introduce LUMA, a unified metric that jointly captures the three key aspects of MU: utility, efficacy, and efficiency. (ii). We validate LUMA on a benchmark of 12 machine unlearning methods on 8 classification datasets spanning four domains (the most comprehensive evaluation of MU methods to date), including the first benchmarks for tabular and graph data. (iii). We propose a taxonomy of the 12 benchmarked methods to aid in structuring and understanding the field. (iv). The benchmark we provide publicly is reproducible, easily extensible, and intended to serve as a foundation for future research in the MU field of study.

Code is available at this anonymized repository, with results reproducible via reproduce.sh.

#### 2 RELATED WORK

As an emerging field, machine unlearning has still seen limited systematic evaluation of the extensive research developed in recent years. In addition, most existing benchmarks remain restricted to a single domain. We list these studies in Table 1, providing details on their coverage (in terms of datasets and methods), the domains they include, and the number of aspects of MU they capture with their metrics.

Several works only focus on one domain: Grimes et al. (2024); Choi & Na (2023); Cadet et al. (2024) evaluate and compare unlearning approaches exclusively on image datasets. Similarly, Koudounas et al. (2025) focuses solely on spoken language understanding datasets.

While valuable, no method was tested on multiple domains jointly, an omission that critically limits our understanding of the generalizability of MU methods. Moreover, most benchmarks either do not consider all the evaluation dimensions or are very limited in the number of datasets and methods they employ.

To the best of our knowledge, the benchmark introduced by Cheng & Amiri (2024) remains the only existing effort to evaluate MU methods across multiple domains. However, it lacks a crucial component of MU evaluation: the incorporation of metrics that directly assess the efficacy of the unlearning process. In their absence, the reported results lack rigorous quantification of unlearning effectiveness, undermining the ability to address the core privacy and security goals that MU is

	Coverage			Modalities					Metrics		
	# datasets	# methods	Image	Text	Tabular	Graph	Speech	Video	Utility	Efficacy	Efficiency
Our	8	12	1	1	1	1	Х	Х	1	1	1
Cheng & Amiri (2024)	6	5	1	_/	Х	Х			/	Х	1
Choi & Na (2023)	2	7	1	X	X	X	X	Х	/	/	×
Cadet et al. (2024)	5	$11+7^{1}$	1	X	X	X	X	Х	/	/	/
Grimes et al. (2024)	1	6	1	X	Х	X	Х	Х	1	Х	×
Koudounas et al. (2025)	5	8	Х	X	×	X	/	Х	1	/	/

Table 1: Comparison of classification benchmarks present in the literature.

meant to achieve (Xu et al., 2024). In addition, as shown in Table 1, their study covers a far narrower set of methods (5 versus our 12), which restricts both the scope and generalizability of their findings.

Table 1 shows that no benchmark in the literature has assessed MU methods in the Tabular and Graph domains. Conversely, our benchmark fills this gap by considering 8 datasets (2 per domain, more than any other benchmark) and 12 methods across four different data domains: Image, Text, Tabular, and Graph. Moreover, we incorporate all key dimensions of MU evaluation (utility, efficacy, and efficiency) within a single framework, introducing LUMA (see Sec. 3.3.1), the first unified multidimensional metric. By providing the most extensive and systematic benchmark to date, we bring coherence to the fragmented landscape of MU in classification tasks and establish an extensible reference point for validating future methods.

#### 3 BENCHMARK DESIGN

We formalize the standard MU workflow as follows: a model M is first trained on a dataset D. Upon receiving an unlearning request, a subset of samples to be removed is specified as the forget set  $D_f$ , with the retain set defined as  $D_r = D \setminus D_f$ . The goal of an unlearning method (Unlearner) is to transform M into an updated model M' that closely approximates a retrained model Gold (Gold Model), i.e., the one obtained by retraining from scratch on  $D_r$ .

We designed our benchmark to faithfully implement this workflow uniformly across datasets. For each method, we perform hyperparameter selection over the learning rate in the range  $10^{-6}$  to  $10^{-3}$ , and additionally tune method-specific parameters where relevant (e.g., the dampening constant in Selective Synaptic Dampening Foster et al. (2024)). All experiments were run three times with different seeds to account for statistical variation.

#### 3.1 Datasets and models

We evaluate all unlearning methods across four domains – image, text, tabular, and graph – using two publicly available datasets per domain. For domains with established MU benchmarks, we adopt datasets widely used in prior work (Chundawat et al., 2023; Golatkar et al., 2020; Foster et al., 2024; Tarun et al., 2023; Fan et al., 2023; Goel et al., 2022; Cha et al., 2024; Cheng & Amiri, 2024). For domains less explored in this context, we select representative and widely studied datasets (Le Quy et al., 2022; Wu et al., 2018).

The forget set  $D_f$  for each dataset is constructed following one of two strategies: (i) selecting training samples containing predefined named entities (e.g., person or organization names) or identity-related attributes (**NE**), simulating realistic deletion requests; or (ii) sampling 20% of the training set uniformly at random (**SA**) when identity-based filtering is not feasible. In both cases, the forget set size is approximately 20% of the training data, consistent with prior MU literature.

For the **image** domain we selected Cifar-100 (Krizhevsky (2009)) (SA) and CelebA (Liu et al. (2015)) (NE). CelebA was used for multilabel classification. For **text** we selected IMDB (Giobergia (2023)) (NE) and AG News (Gulli (2005)) (NE) for **tabular** data, we selected the datasets of Adult (Becker & Kohavi (1996)) (SA) and Spotify Tracks (maharshipandya (2023)) (SA), and for **graphs** we selected BBBP (Sakiyama et al. (2021)) (SA) and BACE (Wu et al. (2018)) (SA). More details on these datasets are reported in Section B.1.

<sup>&</sup>lt;sup>1</sup>7 of the methods reported in this survey were taken from a Machine Unlearning competition and, as such, were not formally peer-reviewed.

For each domain, we adopt model architectures widely used in the MU literature, selecting both smaller and larger variants. In the **image** domain, we use ResNet-18 and ResNet-50 (He et al., 2016), consistent with prior MU studies (Foster et al., 2024; Fan et al., 2023; Golatkar et al., 2020; Tarun et al., 2023). For **text**, we fine-tune DistilBERT (Sanh et al., 2020) and BERT (Devlin et al., 2019) (110M+ parameters), each augmented with a classification head. In the **tabular** domain, we employ two fully connected networks with one and three hidden layers, respectively, each layer containing 100 units. Finally, for **graphs**, we use two GCN-based classifiers: one with a backbone of one GCN layer followed by one dense layer, and another with two layers followed by one dense layer. Full training configurations and hyperparameters are detailed in Appendix B.2.

## 3.2 UNLEARNERS

For this benchmark, we include 12 different state-of-the-art unlearning methods. To facilitate analysis and discussion, we categorize these into four groups based on their unlearning strategy:

Fine Tuning (FT): Gradient Descent (GD), Successive Random Labels (SRL), NegGrad (NG) (Golatkar et al. (2020)), Advanced NegGrad (ANG) (Choi & Na (2023)), UNSIR (UNSIR) (Tarun et al. (2023)). These methods train the model further according to specific strategies.

**Selective Weight Modification (SWM)**: CF-k (CFk) (Goel et al. (2022)), EU-k (EUk) (Goel et al. (2022)), Saliency Unlearning (SalUn) (Fan et al. (2023)). These methods only operate on a subset of the model parameters, typically the last layers or weights identified via saliency masking.

**Distillation** (DIS): Bad Teaching (BT) (Chundawat et al. (2023)), SCRUB (SCRUB) (Kurmanji et al. (2024)). These methods rely on teacher–student setups to alter the target model's behavior.

Weight Importance (WI): Fisher Forgetting (FF) (Golatkar et al. (2020)), Selective Synaptic Dampening (SSD) (Foster et al. (2024)). These directly modify model weights, leveraging the Fisher Information Matrix to estimate the importance of parameters with respect to  $D_f$ .

For reference, the metrics related to the *Original* (Orig.) model and the *Gold Model* (Gold) are also reported for each dataset. In a few cases, certain Unlearners required minor modifications (e.g., change of loss function) to work consistently across settings.

## 3.3 METRICS

We evaluate unlearning methods along three key dimensions:

**Utility**: the predictive performance of the model after unlearning to evaluate the unintended degradation on non-forgotten samples. We compute the F1 score on both the test set and the forget set.

**Efficacy**: the extent to which the influence of the forget set is removed. We adopt the Unlearning Membership Inference Attack (UMIA) (Hayes et al., 2024), which tests whether the model can distinguish forgotten samples from previously unseen ones. Effective unlearning yields UMIA values close to those of the Gold model, avoiding both *over-* and *under-*unlearning (Shi et al., 2024).

**Efficiency**: the computational cost of unlearning. We measure the relative training speedup compared to full retraining, as well as the peak GPU memory usage during execution.

In the following, we refer to Utility, Efficacy, and Efficiency as Evaluation Dimensions (*ED*). While we refer to single benchmarking scores chosen for the *ED*s (e.g., the UMIA for Efficacy) as *measures*. Compound metrics will be referred to simply as *metrics*.

#### 3.3.1 LUMA: A UNIFIED METRIC

Although examining the three *ED*s separately offers a detailed view of an Unlearner's quality, a unified MU metric is essential for comprehensively comparing methods, whether to discard underperforming hyperparameters or to identify the best performer in real-world applications.

Despite its importance, a unified metric remains a critical gap in the MU literature. Koudounas et al. (2025) introduces GUM, a Global Unlearning Metric, but it faces limitations on scalability to multiple measures and resilience to edge cases. GUM reduces each dimension to a single proxy measure (e.g., UMIA for Efficacy, F1 score for Utility, RunTime for Efficiency), which fails to cap-

ture the richness of MU performance when multiple measures per *ED* are available: in practice, MU performance is evaluated by several measures (e.g., F1 scores on forget and test sets, RunTime and GPU memory usage), and considering all these aspects jointly provides a more faithful assessment than relying on single-value proxies. Zhao et al. (2024) proposed Tug of War (ToW), defined as the product of the relative differences between the Gold and retrained models on accuracies over the test, retain, and forget sets. However, ToW excludes Efficiency, leading to a partial evaluation. We provide empirical evidence of these shortcomings in Appendix A.1.

To address these limitations, we introduce the **Laplacian Unlearning Multidimensional Assessment (LUMA)**. Unlike GUM and ToW, LUMA is multidimensional by design, incorporating vectors of measures within each *ED*, and flexible, allowing users to add any measure deemed relevant. LUMA computes the distance of an unlearned model from the retrained Gold model across utility, efficacy, and efficiency, while penalizing large deviations in any single dimension.

All EDs are referenced against the Gold Model (Gold), which represents the ideal target of MU methods (Xu et al. (2024)) as the model retrained from scratch only on the retain set. Let Gold and M' be the Gold Model and the model obtained via the application of the MU method to be measured, respectively.

In our benchmark, the efficacy dimension is measured by UMIA, the utility dimension by F1 score on the test and forget sets, and the Efficiency dimension by runtime and GPU memory usage. Each ED can therefore be represented as a vector of the values computed by the chosen measures on the considered model  $\bullet \in \{\texttt{Gold}, M'\}$ 

$$\mathbf{e}_{\bullet} = [\, \mathrm{UMIA}_{\bullet} \,]^{\top}, \ \mathbf{u}_{\bullet} = \left[ F1^{test}_{\bullet}, \, F1^{forget}_{\bullet} \right]^{\top}, \ \mathbf{t}_{\bullet} = \left[ \mathrm{RunTime}_{\bullet}, \, \mathrm{Memory}_{\bullet} \right]^{\top}$$

To compare the Unlearners against the Gold Model, we map efficacy and utility into similarity factors using a  $\gamma$ -parametrized Laplacian kernel. This choice penalizes *under*- or *over*-performance relative to the Gold Model, while amplifying large deviations in any individual measure.

$$M_U(\mathbf{u}_{\text{Gold}}, \mathbf{u}_{M'}) = \exp((-\gamma \|\mathbf{u}_{\text{Gold}} - \mathbf{u}_{M'}\|_1)),$$
  

$$M_E(\mathbf{e}_{\text{Gold}}, \mathbf{e}_{M'}) = \exp((-\gamma \|\mathbf{e}_{\text{Gold}} - \mathbf{e}_{M'}\|_1)).$$

While their closeness to the Gold Model evaluates utility and efficacy, efficiency follows a different principle: the less time and memory a method consumes, the better; the Gold Model provides an upper bound reference point. To capture this, we first normalize each efficiency measure by the corresponding value of the Gold Model. We then apply a logarithmic transformation to smooth extreme differences and emphasize relative improvements over absolute ones. Finally, we apply an exponential decay so that larger deviations from the Gold baseline are penalized more strongly, and aggregate the resulting values through a weighted average to obtain a single efficiency score:

$$M_T(\mathbf{t}_{\text{Gold}}, \mathbf{t}_{M'}) = \sum_i \mathbf{w}_i \exp \left[ -\left( \frac{\log(1 + \mathbf{t}_{M',i})}{\log(1 + \mathbf{t}_{\text{Gold},i})} \right)^3 \right],$$

where  $\mathbf{w}_i$  represents the weight assigned to the  $i^{th}$  efficiency measure.

Finally, we define the unified metric as:

**Definition 1** (LUMA). Let  $M_U, M_E, M_T$  be the similarity scores for utility, efficacy, and efficiency relative to the Gold model (as defined above). Then

$$LUMA = \frac{3}{\frac{1}{M_U} + \frac{1}{M_E} + \frac{1}{M_T}}.$$

In other words, LUMA is the harmonic mean of  $M_U$ ,  $M_E$ , and  $M_T$ . LUMA is parametrized by the  $\gamma$  of the Laplacian kernels and the weight vector  $\mathbf{w}$  assigned to the efficiency measures. In our setting, we set  $\gamma=3$  and  $\mathbf{w}=(0.9,0.1)$ , assigning 90% of the weight to RunTime and 10% to memory efficiency. We detail parameter tuning of LUMA in Section A.3.

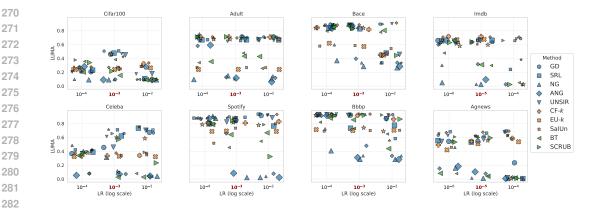


Figure 2: Performance of unlearning methods across datasets under varying learning rates. Each subplot represents a dataset. The marker size denotes the underlying model dimension (small vs. large), while the color indicates the unlearner family (FT, SWM, DIS, WI). The LR used to train the original model is highlighted in dark red.

Note that since the Gold Model is also scored (with a value of 1 for both Utility and Efficacy, as the difference w.r.t. itself is, by definition, 0), its LUMA is fixed. Consequently, any Unlearner that takes longer than the Gold Model will inevitably be penalized by receiving a lower LUMA score, even if it achieves perfect similarity with the Gold Model. This property makes it straightforward to identify cases where an Unlearner is not suitable for application.

LUMA offers several advantages: (i). It is designed to range from 0 to 1, where 1 is the ideal MU algorithm that returns a model identical to the Gold Model with no additional cost in time or memory usage. (ii). It strongly penalizes any significant deviation on any measure, enabling fast pruning of ill-defined (or ill-parameterized) Unlearners. (iii). It is extensible, allowing the integration of any task-specific measure. For example, ROC can replace F1 in the Efficacy ED, or be added alongside it without further modifications. This flexibility is unique to LUMA and ensures seamless integration with future evaluation protocols.

#### 4 RESULTS

In this Section, we summarize the main findings from our benchmark study. In Section 4.1, we start by studying hyperparameters and we show that the size of the model has little impact on the performance of Unlearners. In Section 4.2, we report the main results from our experimentation on the four domains (Tabular, Image, Textual, Graphs). Finally, Section 4.3 summarizes takeaways.

#### 4.1 HYPERPARAMETER TUNING AND MODEL SELECTION

Most of the Unlearners are only parametrized by their Learning Rate (LR) of either the fine-tuning to be applied, the distillation, or any semi-Newton step they employ (refer to Section 3.2 for a taxonomy). In our benchmark, we tuned these parameters by employing LRs one order of magnitude lower, one order of magnitude higher, and equal to the training LR for each domain.

Figure 2 shows the LUMA of each unlearner across datasets and models as the learning rate varies. We use marker size to represent the underlying model dimension. Figure 2 shows that the Negative Gradient unlearners (NG and ANG) are the most sensitive to changes in the LR parameter, often severely degrading performance unless the unlearning LR is set lower than the training LR. In general, using a larger LR than that of training leads to poor LUMA, while the best performance is usually achieved with a rate close to the training value—particularly evident in the Image (CIFAR-100, CelebA) and Text (IMDB, AG News) domains. This makes the tuning of Unlearner dependent on LR quite simple, as the most reliable choice is always to select an LR comparable to the one used for training.

Table 2: Comparison of Unlearners on the Image datasets trained on ResNet-50.

Group	Method	I	Cif	ar 100		I	C	eleba	
Group	Method	F1 test	UMIA	Runtime	LUMA	F1 test	UMIA	Runtime	LUMA
_	Orig.	$.466 \pm .007$	$.706 \pm .028$	$3087.0 \pm 442.3$	$.266 \pm .046$	$.688 \pm .009$	$.770 \pm .011$	$8967.0 \pm 744.0$	$.303 \pm .027$
_	Gold	$.429 \pm .003$	$.501 \pm .003$	$3087.0 \pm 442.3$	$.636 \pm .000$	$.670 \pm .009$	$.539 \pm .008$	$8967.0 \pm 744.0$	$.636 \pm .000$
	GD	$.437 \pm .001$	$.656 \pm .034$	$112.3 \pm 22.2$	$.461 \pm .077$	$.578 \pm .039$	$.535 \pm .006$	$291.3 \pm 31.8$	$.673 \pm .117$
	SRL	$.438 \pm .009$	$.649 \pm .035$	$141.5 \pm 21.8$	$.476 \pm .072$	$.607 \pm .011$	$.529 \pm .005$	$291.8 \pm 32.7$	$.734 \pm .011$
FT	NG	$.001 \pm .000$	$.499 \pm .001$	$38.3 \pm 6.1$	$.092 \pm .004$	$.687 \pm .004$	$.733 \pm .006$	$1.5 \pm .1$	$.397 \pm .023$
	ANG	$.154 \pm .011$	$.536 \pm .016$	$224.1 \pm 32.7$	$.282 \pm .027$	$.216 \pm .055$	$.750 \pm .099$	$693.4 \pm 26.1$	$.071 \pm .019$
	UNSIR	$.436 \pm .005$	$.652 \pm .033$	$317.6 \pm 29.6$	$.452 \pm .071$	$.589 \pm .044$	$.529 \pm .008$	$294.1 \pm 19.2$	$.704 \pm .091$
	CFk	$.462 \pm .009$	$.693 \pm .024$	$89.9 \pm 18.7$	$.322 \pm .029$	$.697 \pm .001$	$.772 \pm .008$	$250.1 \pm 24.6$	$.340 \pm .024$
SWM	${\scriptscriptstyle m EU} k$	$.476 \pm .004$	$.697 \pm .027$	$1733.0 \pm 313.3$	$.264 \pm .031$	$.697 \pm .002$	$.768 \pm .006$	$1269.0 \pm 78.4$	$.329 \pm .013$
	SalUn	$.434 \pm .021$	$.659 \pm .048$	$153.2 \pm 1.6$	$.460 \pm .147$	$.682 \pm .005$	$.620 \pm .020$	$430.3 \pm 199.6$	$.600 \pm .059$
DIS	BT	$.074 \pm .014$	$.651 \pm .013$	$286.7 \pm 5.2$	$.202 \pm .042$	$.318 \pm .091$	$.618 \pm .021$	$744.4 \pm 69.6$	$.340 \pm .128$
DIS	SCRUB	$.460 \pm .025$	$.679 \pm .038$	$143.2 \pm 15.6$	$.339 \pm .100$	$.648 \pm .016$	$.638 \pm .008$	$418.9 \pm 126.4$	$.570 \pm .019$
WI	FF	.186 ± .161	$.065 \pm .922$	$2340.0 \pm 121.8$	$.341 \pm .294$	$.617 \pm .033$	$.655 \pm .024$	$205.6 \pm 26.5$	$.575 \pm .071$
	SSD	$.466 \pm .007$	$.681 \pm .030$	$120.9 \pm .8$	$.316 \pm .066$	$.688 \pm .009$	$.773 \pm .013$	$280.2 \pm 40.3$	$.350 \pm .040$

Table 3: Comparison of Unlearners on the Tabular datasets trained on MLP.

Group	Method	1	A	dult		I	Spo	otify				
Group	Method	F1 test	UMIA	Runtime	LUMA	F1 test	UMIA	Runtime	LUMA			
-	Orig.	$.793 \pm .002$	$.498 \pm .001$	$135.7 \pm 7.9$	$.631 \pm .000$	$.635 \pm .009$	$.528 \pm .005$	$119.6 \pm 1.7$	$.573 \pm .007$			
-	Gold	$.791 \pm .002$	$.499 \pm .001$	$135.7 \pm 7.9$	$.636 \pm .000$	$.629 \pm .002$	$.497 \pm .003$	$119.6 \pm 1.7$	$.636 \pm .000$			
	GD	$.791 \pm .003$	$.498 \pm .002$	$88.3 \pm 4.9$	$.710 \pm .001$	$.604 \pm .007$	$.498 \pm .004$	$4.7 \pm .1$	$.915 \pm .018$			
	SRL	$.790 \pm .002$	$.500 \pm .000$	$11.2 \pm 4.7$	$.668 \pm .002$	$.634 \pm .007$	$.521 \pm .002$	$7.8 \pm .1$	$.842 \pm .012$			
FT	NG	$.434 \pm .000$	$.500 \pm .002$	$16.2 \pm .4$	$.150 \pm .001$	$.604 \pm .016$	$.522 \pm .008$	$1.5 \pm .0$	$.875 \pm .006$			
	ANG	$.767 \pm .004$	$.499 \pm .000$	$140.7 \pm 1.5$	$.591 \pm .006$	$.608 \pm .009$	$.503 \pm .003$	$11.9 \pm .1$	$.875 \pm .002$			
	UNSIR	$.792 \pm .003$	$.498 \pm .001$	$105.2 \pm .6$	$.668 \pm .008$	$.603 \pm .007$	$.499 \pm .001$	$8.5 \pm .2$	$.885 \pm .012$			
	CFk	$.791 \pm .003$	$.499 \pm .000$	$85.5 \pm 1.1$	$.711 \pm .007$	$.642 \pm .007$	$.519 \pm .007$	5.3 ± .1	$.829 \pm .011$			
SWM	EU $k$	$.793 \pm .002$	$.498 \pm .001$	$848.4 \pm 15.6$	$.242 \pm .010$	$.642 \pm .007$	$.514 \pm .007$	$59.6 \pm .3$	$.692 \pm .000$			
	SalUn	$.786 \pm .003$	$.500 \pm .001$	$100.8 \pm 11.9$	$.676 \pm .014$	$.635 \pm .009$	$.526 \pm .006$	$8.3 \pm .6$	$.817 \pm .021$			
DIS	BT	$.660 \pm .106$	$.504 \pm .005$	$156.1 \pm .9$	$.432 \pm .154$	$.583 \pm .027$	$.531 \pm .003$	$13.2 \pm 1.0$	$.785 \pm .029$			
DIS	SCRUB	$.789 \pm .003$	$.498 \pm .002$	$81.4 \pm 1.0$	$.711 \pm .007$	$.611 \pm .013$	$.515 \pm .007$	$7.6 \pm .1$	$.841 \pm .004$			
WI	FF	$.786 \pm .007$	$.500 \pm .001$	$8.0 \pm .8$	$.936 \pm .013$	$.600 \pm .028$	$.517 \pm .008$	$19.4 \pm .2$	$.812 \pm .017$			
vv ⊥	SSD	$.793 \pm .002$	$.499 \pm .001$	$91.5 \pm .9$	$.697 \pm .016$	$.635 \pm .009$	$.518 \pm .007$	$8.1 \pm .1$	$.827 \pm .017$			

The only Unlearners that take as input a parameter different from the LR are the ones in the group of Weight Importance (WI): Fisher Forgetting (FF) and Selective Synaptic Dampening (SSD). In this case, the tuning was applied to their respective parameters, and its results are reported in Section C of the Appendix. In the remainder of the paper, we only report the best results for each Unlearner.

The second key takeaway from Figure 2 is that the size of the model has little effect on the performance of the Unlearners, as they yield close scores across all *ED*s and, as a result, they are close in terms of LUMA. Unlearners perform similarly on models trained on a specific dataset, regardless of their sizes. This result is consistent across all domains.

As a result, in the remainder of the paper, we will only report results from the bigger of the two models. For completeness, we report all other results in Section C of the Appendix.

## 4.2 Main Results

In Tables 2, 3, 4, and 5 we report the main results of our benchmark. All experiments were conducted three times on three different seeds, and the tables report the mean results along with their standard deviation. While LUMA is computed using the full set of measures described in Section 3, for the sake of space we only report one representative measure per *ED* (i.e., test F1 score, UMIA, and RunTime), together with the aggregated LUMA for all experiments. However, full results, including runs for smaller models and different parameters, are reported in Section C. From these tables, we draw intra- (Section 4.2.1) and inter-domain (Section 4.2.2) observations.

## 4.2.1 Intra-domain observations

**Image domain.** Table 2 shows that the best-performing Unlearners in the image domain, according to LUMA, are SRL, GD, Salun, and UNSIR, which achieve comparable results. This holds across

378 379

Table 4: Comparison of the Graph datasets. Each entry reports mean  $\pm$  std.

380
381
382
383
384
385
386
387
300

413 414

421

422

423

429

430

431

BBBP BACE Group Method F1 test UMIA Runtime LUMA UMIA Runtime LUMA F1 test  $0.525 \pm 0.025$  $57.3 \pm 27.8$  $0.564 \pm 0.048$  $0.702 \pm \textbf{0.003}$  $0.501 \pm 0.007$  $55.1 \pm 0.2$  $0.617 \pm \textbf{0.002}$  $0.630 \pm 0.004$ Orig.  $0.559 \pm 0.044$  $0.528 \pm 0.009$  $0.712 \pm 0.007$  $0.498 \pm \textbf{0.005}$  $0.636 \pm 0.000$ Gold  $57.3 \pm 27.8$  $0.636 \pm 0.000$  $55.1 \pm 0.2$  $0.499 \pm 0.006$  $0.612 \pm 0.021$  $0.519 \pm 0.017$  $3.1 \pm 4.0$  $0.705 \pm 0.005$  $1.0 \pm 0.0$  $0.851 \pm 0.102$  $0.924 \pm 0.014$ GD  $0.930 \pm \textbf{0.013}$ SRL  $0.609 \pm 0.015$  $0.536 \pm 0.017$  $4.0 \pm 5.1$  $0.832 \pm \textbf{0.089}$  $0.710 \pm 0.005$  $0.496 \pm 0.004$  $1.4 \pm {\scriptstyle 0.0}$  $0.495 \pm \textbf{0.001}$ NG  $0.385 \pm 0.026$  $0.488 \pm \textbf{0.011}$  $0.2 \pm 0.0$  $0.368 \pm \textbf{0.033}$  $0.485 \pm \textbf{0.046}$  $0.3 \pm {\scriptstyle 0.0}$  $0.408 \pm \textbf{0.126}$  $0.7 \pm {\scriptstyle 0.0}$  $0.869 \pm \textbf{0.031}$ ANG  $0.534 \pm 0.028$  $0.496 \pm 0.017$  $0.688 \pm 0.012$  $0.498 \pm \textbf{0.011}$  $0.9 \pm {\scriptstyle 0.0}$  $0.916 \pm 0.022$  $0.617 \pm \textbf{0.009}$  $0.528 \pm \textbf{0.006}$  $0.850 \pm \textbf{0.097}$  $0.705 \pm 0.005$  $0.494 \pm \textbf{0.006}$  $0.928 \pm \textbf{0.009}$ UNSIR  $1.0 \pm 0.0$  $1.3 \pm 0.0$  $0.522 \pm 0.017$  $0.843 \pm 0.112$  $0.488 \pm 0.002$  $0.920 \pm 0.008$ CFk $0.612 \pm 0.021$  $0.705 \pm 0.005$  $1.1 \pm 0.0$  $3.2 \pm 4.1$  $0.537 \pm \textbf{0.015}$  $42.9 \pm 0.2$ SWM EUk $0.618 \pm 0.110$  $0.709 \pm 0.001$  $34.0 \pm 0.0$  $0.713 \pm 0.005$  $0.616 \pm 0.021$  $0.492 \pm 0.006$  $0.609 \pm 0.015$  $0.536 \pm 0.023$  $0.839 \pm 0.101$  $0.497 \pm 0.006$  $0.919 \pm 0.012$ SalUn  $0.711 \pm 0.007$  $1.3 \pm 0.0$  $1.7 \pm 0.1$  $2.9 \pm 0.0$  $0.570 \pm 0.045$  $0.529 \pm 0.041$  $0.866 \pm 0.011$  $0.494 \pm 0.008$ ВТ  $2.2 \pm 0.1$  $0.659 \pm 0.010$  $0.879 \pm 0.028$ DIS  $0.523 \pm \textbf{0.013}$ SCRUB  $0.612 \pm 0.021$  $1.4 \pm 0.0$  $0.857 \pm 0.126$  $0.717 \pm 0.008$  $0.498 \pm 0.004$  $1.9 \pm 0.0$  $0.920 \pm 0.032$  $0.526 \pm 0.017$  $0.830 \pm 0.108$  $0.505 \pm 0.008$  $0.793 \pm 0.147$  $0.607 \pm 0.020$  $6.9 \pm 0.1$  $0.669 \pm 0.066$ 7.0 + 0.0WΙ SSD  $0.630 \pm 0.004$  $0.527 \pm 0.009$  $1.2 \pm 0.0$  $0.823 \pm 0.099$  $0.702 \pm 0.003$  $0.497 \pm 0.003$  $1.6 \pm 0.1$  $0.935 \pm 0.006$ 

Table 5: Comparison of Unlearners on the Textual datasets trained on BERT.

Group	Method		Ι	MDB			AG	news	
Group	Method	F1 test	UMIA	Runtime	LUMA	F1 test	UMIA	Runtime	LUMA
_	Orig.	$.937 \pm .005$	$.549 \pm .004$	$4914.0 \pm 2965.0$	$.565 \pm .006$	$.880 \pm .012$	$.641 \pm .030$	$5200.0 \pm 54.7$	$.399 \pm .056$
-	Gold	$.939 \pm .006$	$.501 \pm .000$	$4914.0 \pm 2965.0$	$.636 \pm .000$	$.906 \pm .005$	$.525 \pm .024$	$5200.0 \pm 54.7$	$.636 \pm .000$
	GD	$.941 \pm .002$	$.539 \pm .001$	$1329.0 \pm 4.1$	$.667 \pm .045$	$.908 \pm .003$	$.568 \pm .009$	$1655.0 \pm 24.9$	$.544 \pm .063$
	SRL	$.939 \pm .000$	$.542 \pm .003$	$1402.0 \pm .2$	$.661 \pm .043$	$.913 \pm .001$	$.545 \pm .006$	$1748.0 \pm .4$	$.554 \pm .061$
FT	NG	$.408 \pm .090$	$.497 \pm .000$	$75.9 \pm 4.8$	$.051 \pm .031$	$.570 \pm .004$	$.852 \pm .000$	$65.3 \pm .1$	$.158 \pm .028$
	ANG	$.935 \pm .003$	$.502 \pm .010$	$271.0 \pm 8.4$	$.660 \pm .056$	$.893 \pm .002$	$.975 \pm .001$	$3414.0 \pm 8.4$	$.185 \pm .006$
	UNSIR	$.933 \pm .010$	$.538 \pm .004$	$1478.0 \pm \textbf{6.1}$	$.656 \pm .036$	$.910 \pm .002$	$.554 \pm .015$	$1814.0\pm.7$	$.575 \pm .095$
	CFk	$.938 \pm .005$	$.549 \pm .006$	$512.7 \pm 2.9$	$.716 \pm .038$	$.905 \pm .005$	$.555 \pm .017$	$599.1 \pm 5.9$	$.586 \pm .076$
SWM	${\scriptscriptstyle m EU} k$	$.939 \pm .004$	$.551 \pm .002$	$1025.0 \pm 2.5$	$.674 \pm .042$	$.907 \pm .004$	$.556 \pm .014$	$1785.0 \pm 17.2$	$.538 \pm .059$
	SalUn	$.926 \pm .000$	$.540 \pm .000$	$1988.0 \pm 0.0$	$.596 \pm .000$	$.912 \pm .001$	$.552 \pm .006$	$2154.0 \pm 1.8$	$.556 \pm .047$
DIS	BT	$.890 \pm .000$	$.540 \pm .000$	$2427.0 \pm 0.0$	$.656 \pm .000$	$.289 \pm .044$	$.784 \pm .140$	$299.0 \pm 78.3$	$.083 \pm .027$
DIS	SCRUB	$.935 \pm .009$	$.543 \pm .001$	$1904.0 \pm 132.5$	$.638 \pm .055$	$.888 \pm .011$	$.563 \pm .003$	$2258.0 \pm 13.4$	$.594 \pm .087$
WI	FF	$.500 \pm .011$	$.498 \pm .002$	$89.2 \pm 3.3$	$.086 \pm .000$	$.142 \pm .058$	$.854 \pm .007$	.3 ± .0	$.009 \pm .001$
N/ T	SSD	$.941 \pm .000$	$.555 \pm .015$	$1492.0 \pm 83.8$	$.642 \pm .065$	$.879 \pm .010$	$.654 \pm .012$	$1793.0 \pm \textbf{46.3}$	$.440 \pm .065$

both datasets, despite their different setups (CIFAR-100 for multiclass and CelebA for multilabel classification), indicating a degree of stability in image-domain Unlearners. This is not surprising, as the image domain is by far the most extensively studied (see Section 2). However, all these methods substantially increase UMIA relative to Gold (up to .659 and .623, compared to .501 and .539), which results in relatively low LUMA values (below 0.5 and lower than Gold). This highlights that the problem remains unsolved: current Unlearners have yet to match gold-level performance without compromising UMIA.

**Tabular domain.** Table 3 shows that FF is the clear winner on the Adult dataset, followed by GD and SSD. These methods also achieve strong performance on the Spotify dataset. In general, WI methods perform very well in this domain, due to the simpler architecture of MLPs. Both FF and SSD compute a Fisher Information Matrix over the network, which is reasonably accurate on an MLP (Karakida & Osawa (2020)). This is corroborated by the fact that Unlearners in the tabular domain obtain substantially higher LUMA values, reflecting the relative simplicity of this setting. Fine-tuning (FT) methods are also effective, except for NG, which consistently stays behind.

**Graph domain.** FT and DIS Unlearners perform best in the Graph domain, as shown in Table 4. All methods exhibit extremely low running times and relatively minor deviations from the F1 test, as GCNs are notoriously able to generalize better even with fewer samples (Yang et al. (2023)), so they aren't impacted as much by the removal of the Forget Set, and the LUMA scores are very high across the board. This holds across both datasets, although it also highlights the difficulty of completely erasing information from the Gold model itself. Overall, unlearning for graph classification remains underexplored and warrants further investigation.

**Textual domain.** Textual domain is where SWM methods shine, as reported in Table 5. This is because BERT (or LLMs in general) usually fine-tune the last classification layer, while the bulk of knowledge is encoded within the deeper transformer architecture. As such, methods that selectively modify weights (CFk EUk, Salun) remove task-specific information without harming the general representations captured by the transformer layers. Interestingly, WI Unlearners are unable to lever-

age the Fisher Information Matrix to correctly identify which weights to modify on LLMs, as shown by the low LUMA score for FF and SSD on AG news. Lastly, FT and DIS methods achieve good performance (with a few exceptions) but are held back by the non-trivial cost of further training a model of this size in its entirety.

#### 4.2.2 Inter-domain observations

No single Unlearner (or group of Unlearners) dominates across all modalities, which reflects both the early stage of MU research and the ongoing search for more reliable methods. Still, some consistent patterns can be observed across domains.

Fine-tuning is often a safe strategy, so FT Unlearners are generally the most reliable, with the notable exception of NG, which often performs worst in terms of LUMA. Because the Forget Sets we define are heterogeneous (spanning multiple classes), a single epoch of negative gradient updates tends to collapse the model, leading to poor F1 scores on the Test Set. ANG, by contrast, mitigates this issue by also taking the Retain Set into account, and should therefore be always preferred.

Similarly, SWM Unlearners perform consistently well, but they rarely achieve the best performance. Among these, EUk is consistently the worst in terms of LUMA, while Salun is the most reliable across all domains. This aligns with its widespread adoption in the literature, where it often serves as a reference baseline.

On the other hand, both DIS and WI Unlearners show inconsistent performance. The former group performs best with smaller models (e.g., the tabular and graph domains). Still, it fails to scale to large pretrained architectures such as ResNet-50 and BERT, as distilling knowledge from such complex models is considerably more challenging, which is consistent with the literature (Marrie et al. (2024); Fang et al. (2025)). The latter, which directly modify model parameters, are particularly effective given the simpler architecture of MLPs compared to CNNs (e.g., ResNet50) or LLMs (e.g., BERT), where kernels and attention mechanisms may introduce unintended side effects. Among these, SSD is generally more reliable.

#### 4.3 MAIN TAKEAWAYS

The main results of our benchmark can be summarized as follows: (i). Unlearners parameterized by learning rate for a semi-Newton step should adopt a value close to that used during training. Setting the learning rate too high leads to severe degradation of utility, effectively rendering the model unusable. (ii). Scaling up model size does not help: large models do not mitigate unlearning weaknesses. (iii). Unlearning remains fundamentally domain-dependent, as no method succeeds across all domains, although Fine-tuning Unlearners are the most reliable. (iv). A utility-efficacy trade-off is unavoidable with current methods (especially so in the Image domain), showing the field lacks methods that can forget without leaking information. (v). Distillation Unlearning is not scalable: it fails on large pretrained models, while selective weight methods are more reliable for these models. (vi). LUMA exposes hidden weaknesses: prior benchmarks overstated progress by ignoring efficiency and efficacy simultaneously.

## 5 CONCLUSION

In this work, we bring order to the landscape of Machine Unlearning methods for classification by providing: (i). the most comprehensive benchmark to date, covering 4 data modalities, 12 Unlearners, 8 datasets, and 8 models; (ii). the first systematic evaluation on the tabular and graph modalities; (iii). a taxonomy of Unlearner methods in the literature, (iv). a novel unified metric, LUMA, which quantifies performance as a single value to support hyperparameter tuning and unlearner selection in practice; and (v). a publicly available, fully reproducible, and extensible benchmark to facilitate fair comparison in future work.

Our results lead to clear guidelines for the design and deployment of Unlearners, and expose critical shortcomings that can only be revealed through cross-modality analysis. Together, these contributions establish a common ground for evaluating and comparing Machine Unlearning approaches for classification. By ensuring reproducibility and extensibility, our work provides a solid basis for future methods to be rigorously evaluated.

## REFERENCES

- Barry Becker and Ronny Kohavi. Adult. UCI Machine Learning Repository, 1996. DOI: https://doi.org/10.24432/C5XW20.
- Xavier F Cadet, Anastasia Borovykh, Mohammad Malekzadeh, Sara Ahmadi-Abhari, and Hamed
   Haddadi. Deep unlearn: Benchmarking machine unlearning. arXiv preprint arXiv:2410.01276,
   2024.
- Sungmin Cha, Sungjun Cho, Dasol Hwang, Honglak Lee, Taesup Moon, and Moontae Lee. Learning to unlearn: instance-wise unlearning for pre-trained classifiers, 2024. URL https://doi.org/10.1609/aaai.v38i10.28996.
  - Jiali Cheng and Hadi Amiri. Mu-bench: A multitask multimodal benchmark for machine unlearning. *arXiv preprint arXiv:2406.14796*, 2024.
  - Dasol Choi and Dongbin Na. Towards machine unlearning benchmarks: Forgetting the personal identities in facial recognition systems. *arXiv* preprint arXiv:2311.02240, 2023.
  - Vikram S Chundawat, Ayush K Tarun, Murari Mandal, and Mohan Kankanhalli. Can bad teaching induce forgetting? unlearning in deep networks using an incompetent teacher. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 7210–7217, 2023.
  - Kate Crawford. *Atlas of AI: Power, Politics, and the Planetary Costs of Artificial Intelligence*. Yale University Press, 2022.
  - Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2019. URL https://arxiv.org/abs/1810.04805.
  - Chongyu Fan, Jiancheng Liu, Yihua Zhang, Eric Wong, Dennis Wei, and Sijia Liu. Salun: Empowering machine unlearning via gradient-based weight saliency in both image classification and generation. *arXiv* preprint arXiv:2310.12508, 2023.
  - Luyang Fang, Xiaowei Yu, Jiazhang Cai, Yongkai Chen, Shushan Wu, Zhengliang Liu, Zhenyuan Yang, Haoran Lu, Xilin Gong, Yufang Liu, Terry Ma, Wei Ruan, Ali Abbasi, Jing Zhang, Tao Wang, Ehsan Latif, Wei Liu, Wei Zhang, Soheil Kolouri, Xiaoming Zhai, Dajiang Zhu, Wenxuan Zhong, Tianming Liu, and Ping Ma. Knowledge distillation and dataset distillation of large language models: Emerging trends, challenges, and future directions, 2025. URL https://arxiv.org/abs/2504.14772.
  - Jack Foster, Stefan Schoepf, and Alexandra Brintrup. Fast machine unlearning without retraining through selective synaptic dampening. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 12043–12051, 2024.
  - Flavio Giobergia. Imdb-id dataset. https://huggingface.co/datasets/fgiobergia/imdb-id, 2023. Accessed: 2025-05-31.
  - Shashwat Goel, Ameya Prabhu, Amartya Sanyal, Ser-Nam Lim, Philip Torr, and Ponnurangam Kumaraguru. Towards adversarial evaluations for inexact machine unlearning. *arXiv* preprint arXiv:2201.06640, 2022.
  - Aditya Golatkar, Alessandro Achille, and Stefano Soatto. Eternal sunshine of the spotless net: Selective forgetting in deep networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9304–9312, 2020.
- Keltin Grimes, Collin Abidi, Cole Frank, and Shannon Gallagher. Gone but not forgotten: Improved benchmarks for machine unlearning. *arXiv preprint arXiv:2405.19211*, 2024.
  - Michael M Grynbaum and Ryan Mac. The times sues openai and microsoft over ai use of copyrighted work. *The New York Times*, 27, 2023.
- Antonio Gulli. Ag's corpus of news articles. http://www.di.unipi.it/~gulli/AG\_corpus\_of\_news\_articles.html, 2005. Accessed: 2025-05-31.

- Jamie Hayes, Ilia Shumailov, Eleni Triantafillou, Amr Khalifa, and Nicolas Papernot. Inexact unlearning needs more careful evaluations to avoid a false sense of privacy. *arXiv preprint arXiv:2403.01218*, 2024.
  - Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2016.
  - Ryo Karakida and Kazuki Osawa. Understanding approximate fisher information for fast convergence of natural gradient descent in wide neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 808–819, Vancouver, Canada, 2020.
  - Alkis Koudounas, Claudio Savelli, Flavio Giobergia, and Elena Baralis. " alexa, can you forget me?" machine unlearning benchmark in spoken language understanding. *arXiv preprint arXiv:2505.15700*, 2025.
  - Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009. URL https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf.
  - Meghdad Kurmanji, Peter Triantafillou, Jamie Hayes, and Eleni Triantafillou. Towards unbounded machine unlearning. *Advances in neural information processing systems*, 36, 2024.
  - Tai Le Quy, Arjun Roy, Vasileios Iosifidis, Wenbin Zhang, and Eirini Ntoutsi. A survey on datasets for fairness-aware machine learning. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 12(3):e1452, 2022.
  - Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
  - maharshipandya. Spotify tracks dataset. https://huggingface.co/datasets/maharshipandya/spotify-tracks-dataset, 2023.
  - Pratyush Maini, Zhili Feng, Avi Schwarzschild, Zachary C Lipton, and J Zico Kolter. Tofu: A task of fictitious unlearning for llms. *arXiv* preprint arXiv:2401.06121, 2024.
  - Alessandro Mantelero. The eu proposal for a general data protection regulation and the roots of the 'right to be forgotten'. *Computer Law & Security Review*, 29(3):229–235, 2013.
  - Juliette Marrie, Michael Arbel, Julien Mairal, and Diane Larlus. On good practices for task-specific distillation of large pretrained visual models, 2024. URL https://arxiv.org/abs/2402.11305.
  - Hiroshi Sakiyama, Masaki Fukuda, and Yasushi Okuno. Prediction of blood-brain barrier penetration (bbbp) based on molecular descriptors of the free-form and in-blood-form datasets. *Molecules*, 26(24):7428, December 2021. doi: 10.3390/molecules26247428.
  - Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter, 2020. URL https://arxiv.org/abs/1910.01108.
  - Weijia Shi, Jaechan Lee, Yangsibo Huang, Sadhika Malladi, Jieyu Zhao, Ari Holtzman, Daogao Liu, Luke Zettlemoyer, Noah A Smith, and Chiyuan Zhang. Muse: Machine unlearning six-way evaluation for language models. *arXiv preprint arXiv:2407.06460*, 2024.
  - 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.
  - Zhenqin Wu, Bharath Ramsundar, Evan N. Feinberg, Joseph Gomes, Caleb Geniesse, Aneesh S. Pappu, Karl Leswing, and Vijay Pande. Moleculenet: a benchmark for molecular machine learning. *Chemical Science*, 9(2):513–530, 2018. doi: 10.1039/C7SC02664A. URL https://pubs.rsc.org/en/content/articlelanding/2018/sc/c7sc02664a. Open Access.

Heng Xu, Tianqing Zhu, Lefeng Zhang, Wanlei Zhou, and Philip S. Yu. Machine unlearning: A survey, 2024. Chenxiao Yang, Qitian Wu, Jiahua Wang, and Junchi Yan. Graph neural networks are inherently good generalizers: Insights by bridging gnns and mlps, 2023. URL https://arxiv.org/ abs/2212.09034. Kairan Zhao, Meghdad Kurmanji, George-Octavian Barbulescu, Eleni Triantafillou, and Peter Tri-antafillou. What makes unlearning hard and what to do about it. In Advances in Neural Informa-tion Processing Systems (NeurIPS 2024), 2024. 

#### LLM USAGE DISCLOSURE

During the preparation of this work, the authors used LLMs to correct typos and grammatical mistakes. After using this tool/service, the authors reviewed and edited the content as needed and take full responsibility for the content of the published article.

#### A LUMA AGAINST OTHER METRICS

In this section, we show shortcomings and issues of the previous two unified unlearning metrics, GUM(Koudounas et al. (2025)) and ToW (Zhao et al. (2024)). Then, we detail how LUMA handles edge cases and its sensitivity to hyperparameters.

## A.1 SHORTCOMINGS OF OTHER METRICS

Indicating the Gold Model with g and the resulting Unlearned model with u, GUM is defined as:

$$GUM = \frac{(1 + \alpha + \beta)UET}{\alpha ET + \beta UT + UE}$$

$$\begin{aligned} \text{where } U &= 1 - \left| F1_T^{(g)} - F1_T^{(u)} \right|, \quad E &= 1 - \left( \frac{\text{MIA}'(u) - \text{MIA}'(g)}{\text{MIA}(o) - \text{MIA}'(g)} \right)^2, \text{ with} \\ \text{MIA}'(u) &= \min\{\text{MIA}(u), \, \text{MIA}(o)\}, \quad \text{MIA}'(g) = \min\left\{\text{MIA}(g), \, \frac{\text{MIA}'(u) + \text{MIA}(o)}{2} \right\}. \end{aligned}$$

and lastly 
$$T = 1 - \frac{\log(T^{(u)}+1)}{\log(T^{(g)}+1)}$$
.

While GUM includes all the three key evaluation dimensions (*EDs*) for MU methods (efficacy, utility, efficiency) it suffers from two core issues: (i). Only using one measure per dimension: the F1 score on the test set is not enough to assess the model's utility, and the RunTime is not enough to assess the MU method's efficiency. The F1 score on the forget set and the memory efficiency are also important. (ii). When  $MIA'(u) - MIA(o) \le \epsilon$  with a reasonably small  $\epsilon$ , E grows without bound, effectively dominating GUM to unusable values. Worse, the metric is not computable at all for a division by 0 given two conditions:

- $MIA(g) \ge MIA(o)$
- MIA'(u) > MIA(o)

Consider the definition of E, specifically the denominator MIA(o) - MIA'(g). This will be 0 when MIA(o) = MIA'(g). Given the definition of MIA'(g), this can only happen when MIA(g) > MIA(o) (first condition) and MIA'(u) = MIA(o). The latter is true when  $MIA(u) \geq MIA(o)$  (second condition). We show an example in Section A.2.

Tug of War (ToW) is instead defined as:

$$ToW(\theta_u, \theta_g, S, R, D_{test}) = \left(1 - d_a(\theta_u, \theta_g, S)\right) \left(1 - d_a(\theta_u, \theta_g, R)\right) \left(1 - d_a(\theta_u, \theta_g, D_{test})\right),$$
 where

$$a(\theta, D) = \frac{1}{|D|} \sum_{(x,y) \in D} \mathbf{1}[f(x; \theta) = y]$$

is the accuracy of a model f parameterized by  $\theta$  on dataset D, and

$$d_a(\theta_u, \theta_r, D) = |a(\theta_u, D) - a(\theta_r, D)|$$

is the absolute difference between the accuracies of models  $\theta_u$  (the unlearned model) and  $\theta_g$  (the Gold Model) on the dataset D.

ToW includes multiple metrics for the Utility dimension, but misses the Efficiency dimension entirely. Accordingly, the Gold Model will always obtain the optimal ToW score of 1 (regardless of how expensive it is to train) and two unlearners that output the same retrained model with largely varying RunTimes will obtain the same ToW score. Moreover, by omitting MIA, ToW doesn't capture the information leakage of any MU method, which is a known problem in the literature (Xu et al. (2024); Le Quy et al. (2022)).

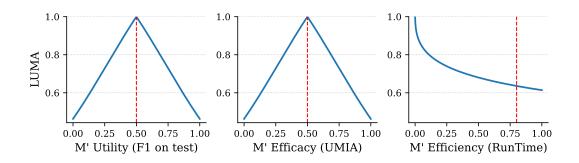


Figure 3: Sensitivity of the proposed LUMA metric with respect to each ED.

#### A.2 EMPIRICAL EVIDENCE

Table 6: Extract of results on the Adult dataset (Table 3).

Group	Method		adult	t3mlp		
Group	Michiga	F1 test	UMIA	Runtime	LUMA	GUM
_	Orig.	$.793 \pm .002$	$.498 \pm .001$	$135.7 \pm 7.9$	$.631 \pm .000$	X
_	Gold	$.791 \pm .002$	$.499 \pm .001$	$135.7 \pm 7.9$	$.636 \pm .000$	X
FT	GD	$.791 \pm .003$	$.498 \pm .002$	$88.3 \pm 4.9$	$.710 \pm .001$	X

Table 6 is an extract from the bigger Table 3 shown in Section 4. This is empirical evidence of the proven problem of GUM in Section A.1: we have that  $MIA(g) \ge MIA(o)$  as  $0.499 \ge 0.498$  (first condition) and that  $MIA'(u) \ge MIA(o)$  as  $0.498 \ge 0.498$ . For this reason, GUM cannot be calculated and will fail because of a statistical variation.

GUM is very useful because it combines measures across all *ED*s, but suffer from numerical instability under some conditions. LUMA fixes this by employing Laplacian kernels.

ToW suffers from the opposite problem: ignoring Efficiency entirely. Consider a toy Unlearner that simply retrains the model from scratch, but deliberately taking double the time to do so. Despite the wasted resources, it would still achieve a perfect ToW score. This illustrates that, to serve as a truly unified metric for unlearning, the *ED*dimension of Efficiency must be incorporated.

#### A.3 LUMA'S BEHAVIOR AND HYPERPARAMETERS

In contrast, LUMA fixes all shortcomings by considering of all the three dimensions, possibly with more than one metric each. By employing Laplacian kernels, LUMA is easily extendable with future MU measures.

In this Section, we analyze the impact of each ED on LUMA, to shed light on its sensitivity and robustness.

By construction, the Laplacian kernels ensure that both positive and negative drifts in an evaluation dimension are penalized symmetrically. Figure 4 illustrates this property, showing the response of LUMA with respect to each *ED* while assuming perfect performance on the remaining *ED*s. For Utility and Efficacy, LUMA reaches its maximum value of 1 when the varying *ED* matches the performance of Gold, and decreases smoothly as the model drifts away in either direction. For Efficiency, LUMA decreases inversely with runtime, lowering as runtime increases.

The Laplacian kernels are parametrized by  $\gamma$ . In this work, we chose  $\gamma=3$  to severely punish even single-measure drifts between the Gold Model (Gold) and the unlearned Model (M').

Figure 4 illustrates the impact of the kernel parameter  $\gamma$  on LUMA as M' drifts from 0 to 1 across all measures. For reference, the Gold model is fixed at 0.5 on all measures. In this work, we chose  $\gamma=3$  as it provided the best smoothness.

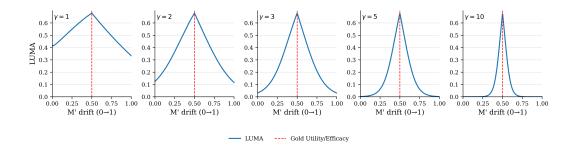


Figure 4: Effect of the Laplacian kernel parameter  $\gamma$  on the LUMA metric as a model drifts from an all-zero representation to an all-one representation of the EDs of Utility and Efficacy.

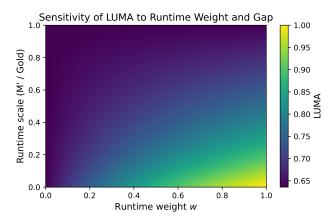


Figure 5: Heatmap of LUMA values as a function of runtime weight w (x-axis) and runtime scale of M' relative to the Gold model (y-axis). The weight assigned to memory is 1-w.

Moreover, LUMA is customizable in terms of weights assigned to measures in the Efficiency ED. The vast majority of works in the Machine Unlearning literature only consider RunTime when evaluating Efficiency, discarding peak memory usage. In this work, we introduced memory usage as a measure with weight = 0.1, to keep RunTime as the most important efficiency metric. This weight vector can be customized according to available resources and the task at hand. Figure 5 shows the sensitivy of LUMA to the weight vector w when the runtime scales w.r.t. to the Gold Model, assuming equality on the Utility and Efficacy EDs. The figure corroborates the correctness of LUMA: if we assign weight 1 to RunTime, and the RunTime is 0, LUMA will be 1, although this is basically unachievable in practice. The vector w, weighting Efficiency measures, is customizable. In this study we set w = [0.9, 0.1], as runtime is the dominant Efficiency factor in unlearning. Figure 5 shows that this choice gives runtime decisive influence over the score while still ensuring memory is not ignored.

Summing up, (i). By employing Laplacian kernels, LUMA fixes the shortcomings of GUM for edge cases, and surpasses both GUM and ToW in considering more than one measure per ED and being easily extensible with other metrics. (ii.) By incorporating all EDs, LUMA surpasses ToW in completeness, as the latter did not consider Efficiency at all.

## B BENCHMARK DETAILS

In this section, we show complete details on the Datasets, Forget Sets, Models, and Unlearners we employed in our benchmark. Moreover, in Section B.3, we show how to easily extend the benchmark for future works.

All details and a step-by-step guide can be found at this anonymized repository.

Dataset	Domain	Features	Samples	Classes
Adult	Tabular	14 tabular features	48,842	2
Spotify	Tabular	15 tabular features	114,000	8
Cifar100	Image	32x32 color images	60,000	100
CelebA	Image	178×218 color images	202,599	40 (ml*)
AG news	Text	Textual descriptions	127,600	4
<b>IMDB</b>	Text	Textual descriptions	100,000	2
BACE	Graph	Graphs ( $\sim$ 34N, $\sim$ 74E, 9F)	1513	2
BBBP	Graph	Graphs ( $\sim$ 23N, $\sim$ 51E, 9F)	2,050	2

Table 7: Overview of the datasets employed for our benchmark. \*ml stands for Multilabel. For Graphs, N stands for nodes, E for edges and F for features.

Architecture	Domain	<b>Epochs</b>	Optimizer	LR
1-hidden layer MLP	Tabular	2	Adam	1e-3
3-hidden layer MLP	Tabular	2	Adam	1e-3
ResNet-18	Image	5	Adam	1e-3
ResNet-50	Image	5	Adam	1e-3
DistilBERT	Text	3	AdamW	2e-5
BERT	Text	3	AdamW	2e-5
DistilBERT	Text	3	AdamW	2e-5
1-layer GCN	Graph	50	RMSprop	1e-3
2-layer GCN	Graph	50	RMSprop	1e-3

Table 8: Overview of the models employed for our benchmark.

#### **B.1** DATASETS

Table 7 shows number of samples, number of attributes, number of classes and domain for each dataset. In the Image domain, we employed CIFAR-100 for multiclass classification and CelebA for binary multilabel classification. As detailed in Section 3, we selected these datasets as the most prominent in the Unlearning literature for each domain.

#### B.2 MODELS

Table 8 shows the general architecture of the model we employed. Our rationale was choosing two models per domain with a similar architecture but different in size. Full configurations can be found at this anonymized repository.

## B.3 EXTENDING THE FRAMEWORK

Our benchmark was built to be easily extendable by researchers and practitioners. New experiments can be defined by just extending the configuration files. A full, step-by-step demonstration on how to implement a new Unlearner method is available at this anonymized repository.

#### C COMPLETE RESULTS

In this Section, we report the tables with the complete results across all experiments.

Table 9: Comparison of Unlearners on ResNet18 trained on Cifar100 (mean  $\pm$  std).

							ai 100 (iiieaii ±	-
Group	Method	LR	LUMA	UMIA		F1 (forget)	RunTime	GPU (MB)
_	Orig.	-		$.651 \pm .019$		$.750 \pm .036$	$1017.0 \pm 289.0$	
_	Gold	-	$0.636 \pm 0.000$	$.499 \pm .001$		$.421 \pm .014$	$1017.0 \pm 289.0$	
	GD	1e-04		$.684 \pm .014$		$.890 \pm .023$	$75.5 \pm 6.4$	$8768 \pm 0$
	GD	1	$0.505 \pm 0.005$			$.642 \pm .016$	$78.4 \pm 9.6$	$8271 \pm 0$
	GD	1e-02	$.286 \pm .067$	$0.502 \pm 0.002$		$.156 \pm .052$	$75.3 \pm 8.3$	$7772 \pm 0$
	SRL	1e-04				$.877 \pm .023$	$110.1 \pm 7.5$	$10256 \pm 0$
	SRL	1e-03	$.508 \pm .039$	$.610 \pm .021$	$1.458 \pm .007$	$.634 \pm .037$	$109.0 \pm 7.2$	$9759 \pm 0$
	SRL	1e-02			$000. \pm 000.$	$000. \pm 000.$	$112.9 \pm 9.2$	$9267 \pm 0$
	NG	1e-04	$.098 \pm .011$	$0.500 \pm 0.001$	$0.001 \pm 0.000$	$.001 \pm .001$	$23.3 \pm 2.6$	$13348 \pm 0$
FΤ	NG	1e-03	$.097 \pm .011$	$0.500 \pm 0.003$	$000. \pm 000.$	$000.\pm000.$	$20.8 \pm 1.0$	$12856 \pm 0$
	NG	1e-02	$.097 \pm .011$	$0.500 \pm 0.000$	$000. \pm 000.$	$000.\pm000.$	$21.9 \pm 1.1$	$12366 \pm 0$
	ANG	1e-04	$.276 \pm .030$	$0.539 \pm 0.024$	$1.151 \pm .018$	$.166 \pm .034$	$150.6 \pm \textbf{24.5}$	$16188 \pm 0$
	ANG	1e-03	$.096 \pm .011$	$0.503 \pm 0.003$	$.001 \pm .001$	$.001\pm.000$	$156.6 \pm 1.5$	$15701 \pm \mathrm{0}$
	ANG	1e-02	$.096 \pm .011$	$0.501 \pm 0.002$	$000. \pm 000.$	$000.\pm000.$	$149.2 \pm \textbf{7.2}$	$15206 \pm 0$
	UNSIR	1e-04	$.205 \pm .008$	$.677 \pm .012$	$0.549 \pm 0.014$	$\textbf{.881} \pm .017$	$183.9 \pm 8.7$	$19149\pm0$
	UNSIR	1e-03	$.469 \pm .029$	$.620 \pm .023$	$0.459 \pm 0.004$	$.649 \pm .036$	$190.1 \pm 4.5$	$18659 \pm 0$
	UNSIR	1e-02	$.177 \pm .054$	$0.502 \pm 0.002$	$0.083 \pm 0.036$	$.084 \pm .036$	$184.1 \pm 7.6$	$18172 \pm \mathrm{0}$
	CF-k	1e-04	$.277 \pm .016$	$.658 \pm .017$	$.520 \pm .006$	$.816 \pm .031$	$76.6 \pm 8.5$	$10498 \pm 0$
	CF-k	1e-03	$.270 \pm .012$	$0.658 \pm 0.016$	$0.524 \pm .007$	$.821 \pm .029$	$82.4 \pm 11.9$	$10228 \pm 0$
	CF-k	1e-02	$.334 \pm .024$	$.645 \pm .016$	$.498 \pm .004$	$.771 \pm .033$	$76.4 \pm 6.7$	$10225 \pm 0$
	EU-k	1e-04		$0.647 \pm 0.015$		$.799 \pm .025$	$752.4 \pm \textbf{26.0}$	$11172\pm0$
SWM	EU- $k$	1e-03	$.248 \pm .008$		$0.525 \pm 0.008$	$.815 \pm .026$	$796.7 \pm 11.4$	$10903 \pm 0$
	EU-k	1e-02		$0.647 \pm 0.016$		$.781 \pm .032$	$794.5 \pm 67.7$	$10768 \pm 0$
	SalUn		$.218 \pm .008$			$.872 \pm .021$	$128.2 \pm 8.6$	$23841 \pm 0$
	SalUn	1	$.426 \pm .062$			$.687 \pm .060$	$145.7 \pm 14.1$	$23257 \pm 0$
	SalUn	1e-02	$.096 \pm .011$	$0.500 \pm 0.000$	$000. \pm 000.$	$000.\pm000.$	$136.0 \pm 4.6$	$22671 \pm 0$
	BT		$0.357 \pm 0.099$	$.570 \pm .005$		$.248 \pm .052$	$190.2 \pm 14.7$	$18055 \pm 0$
	BT	1e-03		$0.569 \pm 0.008$		$.268 \pm .054$	$186.5 \pm 10.0$	$17474 \pm 0$
DIS	BT	1e-02		$0.500 \pm 0.001$		$.002 \pm .002$	$187.3 \pm 11.1$	$16896 \pm 0$
рто	SCRUB	1	$.377 \pm .018$			$.739 \pm .030$	$104.1 \pm 5.0$	$19902 \pm 0$
	SCRUB	1e-03		$0.525 \pm .012$		$.155 \pm .128$	$102.6 \pm 12.4$	$19408 \pm 0$
	SCRUB		$.096 \pm .011$			$000. \pm 000.$	$98.9 \pm 6.9$	$18916 \pm 0$
	FF	1	$1.322 \pm .195$			$.229 \pm .201$	$912.0 \pm 5.4$	$1877 \pm 0$
	FF	1	$.360 \pm .221$	$0.550 \pm .044$		$.258 \pm .224$	$935.9 \pm 44.5$	$1428 \pm 0$
T <sub>N</sub> 7 T	FF	1	$.359 \pm .221$		$1.189 \pm .164$	$.257 \pm .223$	$999.0 \pm 148.4$	$909 \pm 0$
WI	SSD	I	$.357 \pm .022$	$.638 \pm .016$		$.750 \pm .036$	$116.8 \pm 2.5$	$22130\pm0$
	SSD		$.356 \pm .023$			$.750 \pm .036$	$119.4 \pm 6.6$	$21731 \pm 0$
	SSD	1e-02	$.356 \pm .023$	$0.639 \pm 0.016$	$0.489 \pm 0.009$	$.750 \pm .036$	$119.8 \pm 7.2$	$21331 \pm 0$

Table 10: Comparison of Unlearners on ResNet50 trained on Cifar100 (mean  $\pm$  std).

							tar100 (mean =	
Group	Method	LR	LUMA	UMIA		F1 (forget)	RunTime	GPU (MB)
_	Orig.	-	$.266 \pm .046$		$.466 \pm .007$	$.830 \pm .071$	$3087.0 \pm 442.3$	$17695 \pm 0$
_	Gold	-	$0.636 \pm 0.000$		$.429 \pm .003$	$.425 \pm .008$	$3087.0 \pm 442.3$	$17695 \pm 0$
	GD	1e-04	$1.208 \pm .021$		$.509 \pm .008$	$.921 \pm .041$	$108.9 \pm 21.3$	$20924 \pm 0$
	GD	1e-03			$1.437 \pm .001$	$.698 \pm .067$	$112.3 \pm 22.2$	$19880 \pm 0$
	GD	1e-02		$0.500 \pm 0.004$	$1.153 \pm .062$	$.153 \pm .060$	$109.4 \pm 19.8$	$18836 \pm 0$
	SRL	1e-04	$.231 \pm .015$		$.494 \pm .013$	$\textbf{.}897 \pm .035$	$149.3 \pm 29.4$	$24065 \pm \scriptscriptstyle 1$
	SRL	1e-03	$.476 \pm .072$	$.649 \pm .035$	$.438 \pm .009$	$.680 \pm .062$	$141.5 \pm 21.8$	$23018\pm {\scriptscriptstyle 1}$
	SRL	1e-02	$1.09 \pm .028$		$0.022 \pm 0.032$	$.022 \pm .032$	$142.5 \pm 23.8$	$21969 \pm 0$
	NG	1e-04	$0.092 \pm 0.004$	$.499 \pm .001$	$0.001 \pm 0.000$	$.001 \pm .000$	$38.3 \pm 6.1$	$30604 \pm 4$
FT	NG	1e-03	$0.092 \pm 0.004$	$0.500 \pm 0.001$	$000. \pm 000.$	$000. \pm 000$ .	$37.7 \pm 5.6$	$29559 \pm 6$
	NG	1e-02	$.092 \pm .004$	$.500 \pm .001$	$000. \pm 000.$	$000. \pm 000$ .	$37.8 \pm 5.8$	$28513 \pm 6$
	ANG	1e-04	$.282 \pm .027$	$0.536 \pm 0.016$	$.154 \pm .011$	$.180 \pm .027$	$224.1 \pm 32.7$	$39013 \pm 8$
	ANG	1e-03	$0.092 \pm 0.004$	$0.504 \pm 0.004$	$.001 \pm .001$	$.001 \pm .001$	$217.8 \pm 21.9$	$37965 \pm 6$
	ANG	1e-02	$0.091 \pm 0.004$	$0.505 \pm 0.004$	$000. \pm 000.$	$000.\pm000.$	$232.4 \pm 23.7$	$36917\pm 2$
	UNSIR	1e-04	$.215 \pm .014$	$.723 \pm .027$	$0.511 \pm .008$	$.903 \pm .030$	$311.5 \pm 18.8$	$47570 \pm 9$
	UNSIR	1e-03	$.452 \pm .071$	$.652 \pm .033$	$.436 \pm .005$	$.692 \pm .067$	$317.6 \pm 29.6$	$46519 \pm {\scriptstyle 10}$
	UNSIR	1e-02			$0.087 \pm 0.084$	$\textbf{.085} \pm .083$	$320.8 \pm \textbf{35.1}$	$45472 \pm 8$
	CF-k	1e-04		$.711 \pm .030$		$.872 \pm .052$	$84.7 \pm 16.8$	$22545 \pm 1$
	CF-k	1e-03	$0.255 \pm 0.035$	$1.710 \pm .027$	$0.490 \pm .007$	$.874 \pm .051$	$87.1 \pm 19.3$	$21980\pm 2$
	CF-k	1e-02	$.322 \pm .029$		$0.462 \pm 0.009$	$\textbf{.}817 \pm .042$	$89.9 \pm 18.7$	$22358 \pm 6$
	EU- $k$	1e-04			$0.490 \pm .007$	$.866 \pm .060$	$1733.0 \pm 307.7$	$23960 \pm 3$
SWM	EU- $k$	1e-03			$.488 \pm .010$	$.868 \pm .051$	$1730.0 \pm 307.1$	$23395 \pm 4$
	EU- $k$	1e-02			$1.476 \pm .004$	$.845 \pm .051$	$1733.0 \pm 313.3$	$23110 \pm 3$
	SalUn	1e-04		$1.724 \pm .028$		$.892 \pm .030$	$153.4 \pm 1.7$	$11317 \pm {\scriptstyle 1327}$
	SalUn	1e-03		$0.659 \pm 0.048$		$.693 \pm .129$	$153.2 \pm 1.6$	$10081 \pm 1329$
	SalUn	1e-02		$0.500 \pm 0.001$		$000. \pm 000$ .	$153.4 \pm 1.7$	$8844 \pm 1327$
	BT	1e-04			$0.074 \pm .014$	$.172 \pm .054$	$286.7 \pm 50.2$	$40651 \pm 10$
	BT	1e-03			$0.037 \pm .010$	$.129 \pm .034$	$282.2 \pm 43.0$	$39415 \pm 12$
DIS	BT	1e-02			$0.001 \pm 0.001$	$.001 \pm .002$	$274.4 \pm 29.3$	$38177 \pm 11$
סוט	SCRUB	1e-04		$.679 \pm .038$		$.809 \pm .103$	$143.2 \pm 15.6$	$44366 \pm 9$
	SCRUB	1e-03			$1.04 \pm .067$	$.127 \pm .086$	$145.6 \pm 19.7$	$43319 \pm 8$
	SCRUB	1e-02		$.499 \pm .001$	$000. \pm 000.$	$000. \pm 000$ .	$143.2 \pm 17.6$	$42272\pm 8$
	FF	1e-08	$0.341 \pm 0.294$		$.186 \pm .161$	$.282 \pm .244$	$2340.0 \pm 121.8$	$2993 \pm 1429$
	FF	1e-07	$1.337 \pm .291$	$.071 \pm .927$	$.186 \pm .161$	$.290 \pm .251$	$2424.0 \pm 395.7$	$2866 \pm 0$
WI	FF		$.340 \pm .293$		$1.192 \pm .166$	$.297 \pm .257$	$2429.0 \pm 381.4$	$1915 \pm 0$
A A T	SSD		$.316 \pm .066$		$.466 \pm .007$	$.830 \pm .071$	$120.9 \pm .8$	$8118 \pm 500$
	SSD		$.316 \pm .066$		$0.466 \pm 0.007$	$.830 \pm .071$	$121.1 \pm 3.1$	$7254 \pm 499$
	SSD	1e-02	$.314 \pm .064$	$0.682 \pm 0.030$	$0.466 \pm .007$	$.830 \pm .071$	$128.7 \pm 12.5$	$47620 \pm 12$

Table 11: Comparison of Unlearners on ResNet18 trained on CelebA (mean  $\pm$  std).

							elebA (mean ±	
Group	Method	LR	LUMA	UMIA	` ′	F1 (forget)	RunTime	GPU (MB)
_	Orig.	-		$.642 \pm .128$		$.811 \pm .039$	$2202.0 \pm 329.6$	$5379 \pm 2173$
_	Gold	-	$0.636 \pm 0.000$	$.541 \pm .007$	$.589 \pm .064$	$.583 \pm .040$	$2202.0 \pm 329.6$	$5379 \pm 2173$
	GD	1e-04	$1.418 \pm .122$			$.835 \pm .078$	$332.0 \pm 33.1$	$7225 \pm 154$
	GD	1e-03	$0.657 \pm 0.156$	l		$.684 \pm .139$	$352.0 \pm 42.3$	$7018 \pm 281$
	GD	1e-02	$.664 \pm .134$			$.535 \pm .024$	$1453.0 \pm 1826.0$	$6659 \pm 292$
	SRL	1e-04	$.422 \pm .098$	$0.572 \pm 0.086$	$.703 \pm .008$	$.835 \pm .065$	$227.6 \pm 33.8$	$3433 \pm 407$
	SRL	1e-03	$.710 \pm .062$	$0.525 \pm 0.016$	$.638 \pm .041$	$.657 \pm .084$	$229.4 \pm \textbf{26.3}$	$3063 \pm 202$
	SRL	1e-02	$1.707 \pm .118$	$0.529 \pm 0.008$	$.562 \pm .072$	$.548 \pm .035$	$228.1 \pm 41.4$	$2692\pm \imath$
	NG	1e-04	$.485 \pm .154$	$.605 \pm .085$	$.680 \pm .015$	$\textbf{.825} \pm .085$	$1.5 \pm .2$	$7408 \pm 148$
FΤ	NG	1e-03	$.119 \pm .025$	$0.563 \pm 0.023$	$.191 \pm .043$	$.194 \pm .040$	$1.7 \pm .1$	$7276 \pm {\scriptstyle 152}$
	NG	1e-02	$.040 \pm .015$	$.529 \pm .017$	$.045 \pm .022$	$.043 \pm .030$	$4.1 \pm 3.8$	$6905 \pm 152$
	ANG	1e-04	$0.079 \pm 0.032$	$.848 \pm .007$	$.179 \pm .022$	$.122 \pm .012$	$1267.0 \pm \textbf{861.5}$	$8664 \pm 151$
	ANG	1e-03	$0.073 \pm 0.034$	$.779 \pm .017$	$.127 \pm .013$	$.132 \pm .016$	$779.5 \pm 48.9$	$8526 \pm {\scriptstyle 152}$
	ANG	1e-02	$.041 \pm .013$	$0.546 \pm 0.018$	$.052 \pm .012$	$.051 \pm .016$	$1889.0 \pm 1935.0$	$8393 \pm 151$
	UNSIR	1e-04	$.398 \pm .137$	$.585 \pm .107$	$.698 \pm .007$	$\textbf{.}837 \pm .076$	$1084.0 \pm 1271.0$	$10262 \pm 151$
	UNSIR	1e-03	$.614 \pm .183$	$0.553 \pm 0.068$	$.603 \pm .032$	$.683 \pm .117$	$1455.0 \pm 1898.0$	$10129 \pm 153$
	UNSIR	1e-02	$.704 \pm .038$		$.580 \pm .028$	$.561 \pm .021$	$368.7 \pm 29.8$	$9996 \pm 152$
	CF-k		$1.373 \pm .109$		$.701 \pm .011$	$.867 \pm .077$	$338.8 \pm 41.4$	$6742 \pm 202$
	CF-k		$.374 \pm .109$	$.621 \pm .114$		$.870 \pm .074$	$331.7 \pm 4.5$	$6742 \pm 202$
	CF-k	1e-02	$.371 \pm .134$	$.627 \pm .111$	$.694 \pm .008$	$.861 \pm .078$	$1282.0 \pm 1654.0$	$7361 \pm 152$
	$\mathtt{EU-}k$		$1.337 \pm .086$	$.618 \pm .123$	$.700 \pm .011$	$.864 \pm .077$	$1817.0 \pm 37.0$	$6698 \pm 152$
SWM	EU- $k$		$.338 \pm .093$	$.621 \pm .122$	$.700 \pm .011$	$.866 \pm .074$	$1717.0 \pm 203.4$	$6697 \pm 152$
	EU- $k$		$0.349 \pm 0.097$		$.694 \pm .008$	$.859 \pm .080$	$1723.0 \pm \textbf{217.5}$	$6742 \pm 202$
	SalUn		$1.423 \pm .084$		$.703 \pm .009$	$.832 \pm .054$	$210.8 \pm \textbf{26.9}$	$5704 \pm 960$
	SalUn		$1.713 \pm .085$	$0.538 \pm 0.030$		$.639 \pm .059$	$216.1 \pm 37.2$	$5124 \pm 961$
	SalUn		$1.750 \pm .095$			$.603 \pm .028$	$216.8 \pm 28.6$	$4545 \pm 960$
	BT		$0.084 \pm 0.055$			$.141 \pm .119$	$784.9 \pm 71.4$	$8629 \pm 193$
	BT		$0.094 \pm 0.066$	$0.554 \pm 0.006$	$.122 \pm .118$	$.147 \pm .136$	$789.0 \pm 64.9$	$8452 \pm 193$
DIS	BT		$.081 \pm .021$		$.135 \pm .085$	$.150 \pm .094$	$780.4 \pm 82.4$	$8573 \pm 149$
DIO	SCRUB		$.399 \pm .119$		$.707 \pm .011$	$.836 \pm .075$	$442.9 \pm 69.9$	$9218 \pm 148$
	SCRUB		$0.656 \pm 0.008$	$0.545 \pm 0.021$		$.619 \pm .085$	$465.5 \pm 41.2$	$9086 \pm 151$
	SCRUB		$0.386 \pm 0.081$		$.401 \pm .056$	$.369 \pm .026$	$429.7 \pm 66.1$	$8951 \pm 150$
	FF		$0.659 \pm 0.167$	$.554 \pm .054$		$.689 \pm .072$	$76.1 \pm 2.3$	$3308 \pm 1305$
	FF		$0.670 \pm 0.059$		$.576 \pm .029$	$.645 \pm .092$	$75.3 \pm .6$	$2861 \pm 1305$
WI	FF		$.679 \pm .124$	$0.568 \pm 0.052$		$.678 \pm .070$	$75.8 \pm 1.4$	$2416 \pm 1306$
	SSD		$1.419 \pm .140$	l	$.679 \pm .009$	$.842 \pm .090$	$342.7 \pm \textbf{54.0}$	$9277 \pm 147$
	SSD		$.400 \pm .158$			$.842 \pm .090$	$127.0 \pm 1634.0$	$9232 \pm 147$
	SSD	1e-02	$.407 \pm .147$	$.612 \pm .125$	$0.679 \pm 0.009$	$.842 \pm .090$	$685.5 \pm 613.3$	$9291 \pm 146$
		1	•	•				

Table 12: Comparison of Unlearners on ResNet50 trained on CelebA (mean  $\pm$  std).

							DunTime	
Group	Method		LUMA	UMIA		F1 (forget)		GPU (MB)
_	Orig.	-		$.770 \pm .011$	$.688 \pm .009$	$.962 \pm .016$	$8967.0 \pm 744.0$	$15814 \pm 13$
_	Gold	-	$.636 \pm .000$	$.539 \pm .008$	$.670 \pm .009$	$.633 \pm .016$	$8967.0 \pm 744.0$	15814 ± 13
	GD	1e-04			$.704 \pm .003$	$.943 \pm .007$	$284.2 \pm 20.6$	19048 ± 6
	GD			$1.703 \pm .013$		$.892 \pm .006$	$294.4 \pm 37.8$	$18003 \pm 8$
	GD			$0.535 \pm .006$		$.544 \pm .032$	$291.3 \pm 31.8$	$16955 \pm 11$
	SRL			$0.729 \pm 0.010$		$.929 \pm .012$	$401.1 \pm 22.4$	$17758 \pm 7660$
	SRL	1e-03	$.621 \pm .041$	$0.623 \pm 0.028$		$.784 \pm .033$	$302.6 \pm 50.1$	$16709 \pm 7664$
	SRL	1e-02	$.734 \pm .011$	$0.529 \pm 0.005$	$.607 \pm .011$	$.564 \pm .016$	$291.8 \pm 32.7$	$15662 \pm 7669$
	NG	1e-04		$1.733 \pm .006$		$.953 \pm .003$	$1.5 \pm .1$	$27943 \pm 1325$
FT	NG			$0.545 \pm 0.015$		$.421 \pm .055$	$1.6 \pm .1$	$26896 \pm 1323$
	NG	1e-02			$.107 \pm .063$	$.118 \pm .066$	$1.5 \pm .0$	$25850 \pm 1320$
	ANG			$0.878 \pm 0.007$		$.147 \pm .048$	$714.4 \pm 54.1$	$35356 \pm 1329$
	ANG			$1.750 \pm .099$		$.174 \pm .006$	$693.4 \pm 26.1$	$34304 \pm 1328$
	ANG	1e-02	$0.045 \pm 0.018$	$0.558 \pm .013$	$.126 \pm .058$	$.116 \pm .048$	$689.8 \pm 16.6$	$33257 \pm 1328$
	UNSIR	1e-04	$0.369 \pm 0.012$	$0.744 \pm 0.005$	$.705 \pm .003$	$.944 \pm .008$	$296.6 \pm 17.4$	$35105 \pm 12965$
	UNSIR	1		$1.703 \pm .008$		$\textbf{.}878 \pm .058$	$295.2 \pm 18.1$	$34058 \pm 12970$
	UNSIR	1e-02		$0.529 \pm 0.008$		$.564 \pm .061$	$294.1 \pm 19.2$	$33000 \pm 12969$
	CF-k			$.784 \pm .009$		$.969 \pm .001$	$248.1 \pm 20.6$	$20192 \pm 1153$
	CF-k			$1.780 \pm .005$		$.968 \pm .001$	$246.4 \pm 17.5$	$19629 \pm 1153$
	CF-k	1		$1.772 \pm .008$		$.969 \pm .004$	$250.1 \pm 24.6$	$19913 \pm 1658$
	EU-k	1e-04	$0.325 \pm 0.016$	$1.775 \pm .003$	$.698 \pm .001$	$.968 \pm .001$	$1284.0 \pm 68.4$	$22261 \pm 1149$
SWM	EU-k	1e-03	$0.329 \pm 0.013$		$.697 \pm .002$	$.968 \pm .001$	$1269.0 \pm 78.4$	$21696 \pm 1149$
	EU-k	1e-02		$1.780 \pm .008$		$.968 \pm .002$	$1234.0 \pm 93.1$	$20756 \pm 1152$
	SalUn	1e-04		$1.736 \pm .008$		$.932 \pm .010$	$486.0 \pm 167.0$	$15095 \pm 8763$
	SalUn			$0.620 \pm 0.020$		$.806 \pm .023$	$430.3 \pm 199.6$	$27398 \pm 18090$
	SalUn			$0.529 \pm 0.008$		$.511 \pm .027$	$419.5 \pm 216.1$	$26158 \pm 18089$
	BT		$.340 \pm .128$		$.318 \pm .091$	$.607 \pm .150$	$744.4 \pm 69.6$	$30182 \pm 12965$
	BT	1e-03	$0.027 \pm 0.006$	$0.567 \pm .033$		$.102 \pm .082$	$735.1 \pm 66.2$	$28946 \pm 12968$
DIS	BT	1e-02	$0.030 \pm 0.007$		$.074 \pm .036$	$.081 \pm .047$	$711.2 \pm 30.0$	$27715 \pm 12965$
DID	SCRUB	1		$1.751 \pm .005$		$.942 \pm .010$	$371.8 \pm 41.9$	$33746 \pm 12968$
	SCRUB			$0.638 \pm 0.008$		$.808 \pm .033$	$418.9 \pm 126.4$	$32701 \pm 12970$
	SCRUB	1e-02		$0.551 \pm .024$		$.352 \pm .029$	$371.7 \pm 46.2$	$31655 \pm 12969$
	FF	1		$.672 \pm .023$		$.793 \pm .078$	$211.2 \pm 30.2$	$3798 \pm 0$
	FF	1e-07	$0.575 \pm 0.071$	$0.655 \pm 0.024$		$.775 \pm .067$	$205.6 \pm 26.5$	$2841 \pm 0$
WI	FF	1e-06	$0.550 \pm 0.041$	$0.659 \pm 0.020$	$.629 \pm .012$	$.811 \pm .029$	$196.7 \pm 30.7$	$1885 \pm 0$
	SSD	1e-04		$1.773 \pm .013$		$.962 \pm .017$	$280.2 \pm 40.3$	$36018 \pm 12967$
	SSD			$.780 \pm .022$		$.962 \pm .017$	$281.9 \pm 43.9$	$35161 \pm 12967$
	SSD	1e-02	$.348 \pm .044$	$.777 \pm .021$	$.688 \pm .009$	$.962 \pm .017$	$273.8 \pm 28.8$	$34306 \pm 12968$

Table 13: Comparison of Unlearners on DistilBERT trained on AG News (mean  $\pm$  std).

- Orig	Croup Method   I P   I I I I I I I I I I I I I I I I									
- Gold636±.000	Group		LR	LUMA	UMIA			RunTime	GPU (MB)	
GD	_	_								
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	_									
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$										
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$										
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$										
SRL   1e-04   .694 ± .034   .560 ± .013   .924 ± .002   .792 ± .057   557.2 ± 1.3   35357 ± .0669 ± .177   .830 ± .112   .574 ± .156   20.7 ± .4   14608 ± .080   1e-05   .014 ± .005   .857 ± .001   .269 ± .011   .140 ± .080   21.0 ± .0   38420 ± .080   .010 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .21.0 ± .0   37654 ± .080   .010 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000   .011 ± .000										
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$			l						$36124 \pm 0$	
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	FT								$35357 \pm 0$	
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$										
ANG le-06 $.153 \pm .022$ $.975 \pm .002$ $.918 \pm .001$ $.304 \pm .084$ $1087.0 \pm 11.5$ $17788 \pm .084$ ANG le-05 $.079 \pm .001$ $.988 \pm .000$ $.926 \pm .005$ $.014 \pm .006$ $1071.0 \pm 4.3$ $20877 \pm .084$ $.0081 \pm .006$ $.025 \pm .034$ $.071 \pm .026$ $.412 \pm .393$ $.008 \pm .003$ $1088.0 \pm 3.5$ $42367 \pm .084$ $.0851 \pm .096$ $.0627 \pm .017$ $.560 \pm .003$ $.935 \pm .001$ $.954 \pm .007$ $584.1 \pm 6.5$ $16116 \pm 44$ $.0851 \pm .096$ $.080 \pm .093$ $.935 \pm .001$ $.961 \pm .005$ $.080 \pm .093$ $.093 \pm .093$		NG							$38420 \pm 0$	
ANG 1e-05 .079 ± .001 .988 ± .000 .926 ± .005 .014 ± .006		NG							$37654\pm0$	
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$		ANG				$.918 \pm .001$			$17788 \pm 0$	
UNSIR 1e-06 $.627 \pm .017$ $.560 \pm .003$ $.935 \pm .001$ $.954 \pm .007$ $584.1 \pm 6.5$ $16116 \pm 44$ UNSIR 1e-05 $.617 \pm .007$ $.561 \pm .002$ $.931 \pm .001$ $.961 \pm .005$ $579.1 \pm 1.5$ $23294 \pm .008$ UNSIR 1e-04 $.683 \pm .034$ $.559 \pm .012$ $.917 \pm .006$ $.782 \pm .060$ $579.3 \pm 1.7$ $22526 \pm .026$ $CF-k$ 1e-05 $.670 \pm .014$ $.574 \pm .003$ $.935 \pm .002$ $.962 \pm .003$ $211.7 \pm 6.2$ $13089 \pm .026$ $CF-k$ 1e-05 $.675 \pm .022$ $.571 \pm .003$ $.935 \pm .002$ $.957 \pm .011$ $205.0 \pm 1.0$ $34862 \pm .026$ $EU-k$ 1e-06 $.458 \pm .052$ $.748 \pm .047$ $.932 \pm .004$ $.955 \pm .020$ $636.5 \pm 19.9$ $12581 \pm .026$ $EU-k$ 1e-05 $.612 \pm .019$ $.568 \pm .001$ $.935 \pm .002$ $.955 \pm .013$ $615.3 \pm 2.1$ $34862 \pm .026$ $EU-k$ 1e-04 $.626 \pm .019$ $.557 \pm .005$ $.935 \pm .002$ $.949 \pm .010$ $616.2 \pm 1.8$ $34862 \pm .026$ $Salun$ 1e-05 $.619 \pm .032$ $.548 \pm .007$ $.931 \pm .001$ $.940 \pm .038$ $764.2 \pm 23.5$ $32588 \pm .026$ $Salun$ 1e-04 $.634 \pm .041$ $.549 \pm .007$ $.926 \pm .007$ $.733 \pm .045$ $767.6 \pm 13.1$ $31312 \pm .026$ $BT$ 1e-05 $.043 \pm .011$ $.975 \pm .003$ $.123 \pm .040$ $.611 \pm .036$ $972.5 \pm 2.7$ $23190 \pm .028$ $.043 \pm .011$ $.975 \pm .003$ $.123 \pm .040$ $.611 \pm .036$ $.072.5 \pm 2.7$ $.03190 \pm .028$		ANG	1e-05	$.079 \pm .001$	$0.988 \pm 0.000$	$0.926 \pm 0.005$	$.014 \pm .006$	$1071.0 \pm 4.3$	$20877 \pm 0$	
UNSIR 1e-05 $.617 \pm .007$ $.561 \pm .002$ $.931 \pm .001$ $.961 \pm .005$ $579.1 \pm 1.5$ $23294 \pm .008$ $1e-04$ $.683 \pm .034$ $.559 \pm .012$ $.917 \pm .006$ $.782 \pm .060$ $579.3 \pm 1.7$ $22526 \pm .028$ $1e-05$ $.670 \pm .014$ $.574 \pm .003$ $.935 \pm .002$ $.962 \pm .003$ $211.7 \pm 6.2$ $13089 \pm .028$ $1e-05$ $1$			1e-04	$.025 \pm .034$	$.971 \pm .026$	$.412 \pm .393$	$.008 \pm .003$	$1088.0 \pm 3.5$	$42367 \pm 0$	
UNSIR $  1e-04  $ .683 ± .034 .559 ± .012 .917 ± .006 .782 ± .060 .579.3 ± 1.7 .22526 ± .026 .67-k .026 .670 ± .014 .574 ± .003 .935 ± .002 .962 ± .003 .211.7 ± 6.2 .13089 ± .027 .027 .025.0 ± 1.0 .034862 ± .027 .025.0 ± 1.0 .034862 ± .027 .025.0 ± 1.0 .034862 ± .027 .025.0 ± 1.0 .034862 ± .027 .025.0 ± 1.0 .034862 ± .027 .025.0 ± 1.0 .034862 ± .027 .025.0 ± 1.0 .034862 ± .027 .025.0 ± 1.0 .034862 ± .027 .025.0 ± 1.0 .034862 ± .027 .025.0 ± 1.0 .034862 ± .027 .025.0 ± 1.0 .034862 ± .027 .025.0 ± 1.0 .034862 ± .027 .025.0 ± 1.0 .034862 ± .027 .025.0 ± 1.0 .034862 ± .027 .025.0 ± 1.0 .034862 ± .027 .025.0 ± 1.0 .034862 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± 1.0 .034862 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027 .025.0 ± .027		UNSIR				$.935 \pm .001$		$584.1 \pm 6.5$	$16116 \pm 4429$	
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$		UNSIR	1e-05	$.617 \pm .007$	$0.561 \pm 0.002$	$.931 \pm .001$			$23294 \pm 0$	
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$		UNSIR	1e-04	$.683 \pm .034$	$0.559 \pm 0.012$	$.917 \pm .006$		$579.3 \pm 1.7$	$22526 \pm 0$	
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$		CF-k	1e-06	$.670 \pm .014$	$.574 \pm .003$	$.935 \pm .002$			$13089 \pm 0$	
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$		$\mathtt{CF-}k$	1e-05	$.675 \pm .022$	$1.571 \pm .003$	$0.935 \pm 0.002$	$.957 \pm .011$		$34862 \pm 0$	
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$		$\mathtt{CF-}k$	l	$.690 \pm .025$					$34862 \pm 0$	
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$									$12581 \pm 0$	
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	SWM								$34862 \pm 0$	
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$									$34862 \pm 0$	
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$		SalUn							$17408 \pm 4430$	
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$		SalUn	1e-05	$.619 \pm .032$	$0.548 \pm 0.007$	$0.931 \pm 0.001$			$32588 \pm 590$	
BT $  1e-05   .043 \pm .011   .975 \pm .003   .123 \pm .040   .611 \pm .036   972.5 \pm 2.7   23190 \pm .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   .011   $									$31312 \pm 590$	
									$14990 \pm 4430$	
1 04 010   010   100   000   001 0   001 0									$23190 \pm 0$	
	DIG	BT					$.292 \pm .083$	$971.2 \pm \textbf{2.8}$	$22170\pm 0$	
$  SCRUB   1e-00   .008 \pm .014   .301 \pm .003   .930 \pm .001   .930 \pm .003   /31./ \pm 6.4   1343/ \pm 44$	סוט	SCRUB	l						$15457 \pm 4430$	
		SCRUB	l			$0.934 \pm 0.001$	$.963 \pm .002$	$726.8 \pm 2.0$	$25186 \pm 0$	
		SCRUB	1e-04	$.581 \pm .082$	$1.573 \pm .055$	$.869 \pm .067$		$726.7 \pm 1.8$	$24422 \pm 0$	
		FF	1e-08	$.010 \pm .004$					$7220 \pm 0$	
		FF			$0.699 \pm 0.029$	$1.147 \pm .028$	$.111 \pm .111$		$6195 \pm 0$	
	Ta7 T	FF	1e-06	$.005 \pm .001$	$1.711 \pm .158$	$1.136 \pm .032$		$1.0 \pm 1.3$	$5169 \pm 0$	
SSD $ 1e-06 $ $.612 \pm .012 $ $.577 \pm .003 $ $.935 \pm .002 $ $.962 \pm .003 $ $553.4 \pm 4.9 $ $16072 \pm 44$	AA T	SSD	1e-06	$.612 \pm .012$				$553.4 \pm 4.9$	$16072 \pm 4430$	
SSD $ 1e-05 .607 \pm .007 .580 \pm .004 .935 \pm .002 .962 \pm .003 .547.7 \pm 1.4 .29382 \pm .003 .004 .004 .004 .004 .004 .004 .004 $		SSD	1e-05	$.607 \pm .007$	$0.580 \pm 0.004$	$0.935 \pm .002$	$.962 \pm .003$	$547.7 \pm 1.4$	$29382 \pm \mathrm{0}$	
SSD $ 1e-04 $ $.611 \pm .008 $ $.575 \pm .003 $ $.935 \pm .002 $ $.962 \pm .003 $ $548.9 \pm 2.0 $ $29127 \pm .003 $		SSD	1e-04	$.611 \pm .008$	$3.575 \pm .003$	$0.935 \pm .002$	$.962 \pm .003$	$548.9 \pm 2.0$	$29127 \pm 0$	

Table 14: Comparison of Unlearners on BERT trained on AG News (mean  $\pm$  std).

Group	Method	LR	LUMA	UMIA	F1 (test)	F1 (forget)	RunTime	GPU (MB)
_	Orig.	-	$.399 \pm .056$		$.880 \pm .012$		$520.0 \pm 54.7$	$56962 \pm 27786$
_	Gold	-	$.636 \pm .000$	$.525 \pm .024$	$.906 \pm .005$	$.666 \pm .036$	$520.0 \pm 54.7$	$56962 \pm 27786$
	GD	1e-06	$.543 \pm .069$	$.559 \pm .007$	$.913 \pm .001$	$.876 \pm .017$	$1685.0 \pm 1.7$	$32123 \pm 0$
	GD	1e-05	$0.544 \pm 0.063$	$\textbf{.}568 \pm .009$	$0.908 \pm 0.003$	$\textbf{.}867 \pm .019$	$1655.0 \pm 24.9$	$75876 \pm 2154$
	GD	1e-04	$.235 \pm .377$	$.645 \pm .162$	$.436 \pm .412$	$\textbf{.}236 \pm \textbf{.}318$	$1655.0 \pm 25.0$	$75473 \pm 2819$
	SRL	1e-06	$.554 \pm .061$	$\textbf{.545} \pm .006$	$.913 \pm .001$	$\textbf{.866} \pm .018$	$1748.0 \pm .4$	$25943 \pm 0$
	SRL	1e-05	$0.532 \pm 0.000$	$\textbf{.}567 \pm .000$	$0.909 \pm 0.000$	$\textbf{.}872 \pm .000$	$1702.0 \pm .0$	$77502 \pm \mathrm{0}$
	SRL		$000. \pm 800.$		$1.03 \pm .000$	$\textbf{.001} \pm .000$	$1702.0 \pm .0$	$76249 \pm 0$
	NG	1e-06	$1.158 \pm .028$	$\textbf{.}852 \pm .000$	$0.570 \pm .004$	$\textbf{.369} \pm .040$	$65.3 \pm .1$	$27199 \pm 0$
FΤ	NG	1e-05	$0.012 \pm 0.004$	$\textbf{.859} \pm .001$	$1.185 \pm .061$	$\textbf{.011} \pm .000$	$65.6 \pm .1$	$44945 \pm 6656$
	NG	1e-04	$0.009 \pm 0.001$	$.859 \pm .001$	$1.100 \pm .000$	$.011 \pm .000$	$65.6 \pm .0$	$43691 \pm 6657$
	ANG	1e-06	$1.185 \pm .006$	$.975 \pm .001$	$0.893 \pm 0.002$	$.316 \pm .010$	$3414.0 \pm 8.4$	$33283 \pm 0$
	ANG	1e-05	$1.105 \pm .006$	$.986 \pm .001$	$0.867 \pm 0.006$	$\textbf{.011} \pm .000$	$3413.0 \pm 2.4$	$50323 \pm 3865$
	ANG	1e-04	$.009 \pm .001$			$\textbf{.007} \pm .006$	$3413.0 \pm 2.3$	$49065 \pm 3865$
	UNSIR	1e-06	$0.543 \pm 0.072$	$\textbf{.558} \pm .009$	$0.913 \pm 0.000$	$\textbf{.}871 \pm .023$	$1844.0 \pm 51.3$	$34828 \pm 0$
	UNSIR	1e-05	$.575 \pm .095$	$\textbf{.}554 \pm .015$	$.910 \pm .002$	$\textbf{.}832 \pm .050$	$1814.0 \pm .7$	$54186 \pm 3700$
	UNSIR	1e-04	$0.022 \pm 0.022$	$\textbf{.828} \pm .026$	$1.187 \pm .145$	$.080 \pm .137$	$1818.0 \pm 6.3$	$52931 \pm 3700$
	CF-k	1e-06			$.893 \pm .005$	$.902 \pm .027$	$596.3 \pm .4$	$23755 \pm 0$
	CF-k		$0.543 \pm 0.065$			$\textbf{.}895 \pm .016$	$595.7 \pm 5.6$	$74400 \pm 4542$
	CF-k	1e-04	$0.586 \pm 0.076$	$\textbf{.555} \pm .017$	$0.905 \pm 0.005$	$\textbf{.}873 \pm .013$	$599.1 \pm 5.9$	$74400 \pm 4542$
	EU- $k$	1e-06	$.421 \pm .101$	$.723 \pm .051$	$0.889 \pm 0.006$	$\textbf{.}879 \pm .060$	$1789.0 \pm .6$	$22922\pm 0$
SWM	EU- $k$	1e-05	$0.510 \pm 0.059$	$\textbf{.}587 \pm .023$	$0.898 \pm 0.008$	$\textbf{.}882 \pm .036$	$1784.0 \pm 16.7$	$74400 \pm 4542$
	EU- $k$	1e-04	$0.538 \pm 0.059$	$\textbf{.}556 \pm .014$	$0.907 \pm .004$	$\textbf{.}876 \pm .007$	$1785.0 \pm 17.2$	$74400 \pm 4542$
	SalUn	1e-06			$0.913 \pm 0.000$	$.864 \pm .016$	$218.0 \pm 73.0$	$35988 \pm 0$
	SalUn	1e-05	$0.556 \pm .047$	$\textbf{.}552 \pm .006$	$0.912 \pm .001$	$\textbf{.846} \pm .002$	$2154.0 \pm 1.8$	$69079 \pm 6163$
	SalUn		$0.008 \pm 0.001$			$\textbf{.007} \pm .005$	$2152.0 \pm 7.7$	$66978 \pm 6165$
	BT		$.083 \pm .027$		$.289 \pm .044$	$.421 \pm .088$	$299.0 \pm 78.3$	$32106 \pm 0$
	BT		$0.041 \pm 0.009$		$1.104 \pm .002$	$.433 \pm .077$	$2942.0 \pm 2.9$	$57465 \pm 5999$
DIS	BT		$0.017 \pm .009$		$0.097 \pm .004$	$.166 \pm .136$	$2944.0 \pm 4.9$	$55790 \pm 5998$
DIS	SCRUB	1e-06			$0.912 \pm .000$	$.871 \pm .022$	$2287.0 \pm 61.6$	$32623 \pm 0$
	SCRUB	1e-05	$0.594 \pm 0.087$		$1.888 \pm .011$	$.779 \pm .073$		$61258 \pm 6666$
	SCRUB		$0.012 \pm 0.007$		$1.160 \pm .099$	$.004 \pm .004$		$60007 \pm 6666$
	FF		$0.009 \pm 0.001$		$1.142 \pm .058$	$.011 \pm .001$	$.3 \pm .0$	$76257 \pm 2091$
	FF		$0.009 \pm 0.002$			$.020 \pm .027$	$.3 \pm .0$	$75371 \pm 3213$
WI	FF		$100. \pm 800$ .		$1.101 \pm .001$	$.010 \pm .001$	$2.3 \pm .5$	$78113 \pm 4$
AA T	SSD		$.440 \pm .065$		$0.879 \pm 0.010$	$.907 \pm .021$	$1793.0 \pm 46.3$	$33729 \pm 0$
	SSD	1e-05			$0.879 \pm 0.010$	$.907 \pm .021$	$1765.0 \pm 1.1$	$58845 \pm 2894$
	SSD	1e-04	$.439 \pm .063$	$.653 \pm .010$	$0.879 \pm 0.010$	$.907 \pm .021$	$1765.0 \pm 1.0$	$58428 \pm 2894$
							1	

Table 15: Comparison of Unlearners on BERT trained on DistilBERT (mean  $\pm$  std).

	Method		LUMA	UMIA		F1 (forget)	RunTime	GPU (MB)
-	Orig.	-	$.564 \pm .011$	$.548 \pm .002$		$.981 \pm .006$	$949.0 \pm 7.8$	$32839 \pm 0$
_	Gold	-		$.500 \pm .004$			$949.0 \pm 7.8$	$32839 \pm 0$
	GD	1e-06		$.544 \pm .004$		$.988 \pm .002$	$464.9 \pm 3.3$	$38165 \pm 0$
	GD	1e-05		$.545 \pm .006$			$465.1 \pm 3.1$	$34185 \pm 0$
	GD	1e-04		$.512 \pm .006$		$.899 \pm .008$	$465.0 \pm 3.0$	$33383 \pm 0$
	SRL	1e-06		$.545 \pm .002$		$.986 \pm .002$	$490.9 \pm 3.2$	$34978 \pm 0$
	SRL	1e-05		$.544 \pm .003$		$.984 \pm .004$	$491.3 \pm 3.1$	$36570 \pm 0$
	SRL	1e-04		$.508 \pm .016$		$.531 \pm .343$	$490.2 \pm 3.0$	$35775 \pm 0$
	NG	1e-06		$.519 \pm .031$		$.586 \pm .334$	$26.0 \pm .4$	$37371 \pm 0$
FT	NG	1e-05		$.502 \pm .005$			$26.0 \pm .4$	$38962 \pm 0$
	NG	1e-04		$.501 \pm .004$			$25.9 \pm .4$	$38168 \pm 0$
	ANG	1e-06		$.528 \pm .002$			$932.5 \pm 7.3$	$40966 \pm 0$
	ANG	1e-05	$.316 \pm .109$	$.511 \pm .010$	$.828 \pm .081$	$.573 \pm .073$	$928.4 \pm 5.8$	$42569 \pm 0$
	ANG	1e-04	$.026 \pm .001$	$.509 \pm .009$	$0.332 \pm 0.000$	$.331 \pm .006$	$929.2 \pm 5.6$	$41768 \pm 0$
	UNSIR	1e-06	$.628 \pm .011$	$.544 \pm .003$	$.938 \pm .001$	$\textbf{.987} \pm .001$	$519.3 \pm 5.0$	$18614 \pm 0$
	UNSIR	1e-05	$.635 \pm .014$	$.540 \pm .004$	$.936 \pm .002$	$\textbf{.984} \pm .003$	$520.9 \pm 4.8$	$20215\pm \imath$
	UNSIR	1e-04	$.680 \pm .013$	$.518 \pm .005$	$.899 \pm .003$		$521.6 \pm \textbf{5.3}$	$19415\pm 0$
	CF-k	1e-06	$.706 \pm .019$	$.553 \pm .002$	$.934 \pm .002$		$182.3 \pm 2.3$	$36989 \pm 0$
	CF-k	1e-05	$.707 \pm .017$	$.550 \pm .004$	$.935 \pm .001$	$\textbf{.983} \pm .004$	$183.2 \pm 2.2$	$36467 \pm 0$
	CF-k	1e-04	$.706 \pm .014$	$.549 \pm .001$	$.935 \pm .002$	$\textbf{.984} \pm .003$	$182.9 \pm \textbf{2.3}$	$36467 \pm 0$
	EU- $k$	1e-06		$.545 \pm .010$		$\textbf{.981} \pm .006$	$365.2 \pm 4.4$	$36467 \pm 0$
SWM	EU-k	1e-05		$.550 \pm .004$		$\textbf{.983} \pm .003$	$365.3 \pm 4.8$	$36468 \pm 0$
	EU- $k$	1e-04		$.550 \pm .004$		$\textbf{.984} \pm .003$	$365.0 \pm 4.6$	$36467 \pm 0$
	SalUn	1e-06		$.545 \pm .002$		$\textbf{.985} \pm .002$		$14472 \pm 14372$
	SalUn	1e-05		$.542 \pm .006$		$\textbf{.982} \pm .002$	$692.7 \pm 9.5$	$32989 \pm 600$
	SalUn	1e-04		$.505 \pm .003$		$\textbf{.688} \pm .296$	$718.3 \pm \textbf{2.6}$	$31682 \pm 601$
	BT	1e-06		$.508 \pm .009$		$.825 \pm .156$	$879.3 \pm 8.7$	$20328 \pm 0$
	BT	1e-05		$.506 \pm .011$		$.716 \pm .247$	$880.4 \pm 8.6$	$22442\pm 7$
DIS	BT	1e-04		$.504 \pm .002$		$.432 \pm .095$	$876.9 \pm 9.6$	$21387 \pm \mathtt{5}$
DIO	SCRUB	1e-06		$.547 \pm .002$		$.987 \pm .002$	$63.2 \pm 6.0$	$22815\pm \imath$
	SCRUB	1e-05		$.544 \pm .001$		$.975 \pm .012$	$63.5 \pm 6.2$	$24405 \pm 4$
	SCRUB	1e-04		$.512 \pm .006$		$.908 \pm .005$	$63.9 \pm 6.1$	$23608 \pm 3$
	FF	1e-08		$.506 \pm .006$		$.395 \pm .065$	$.1 \pm .0$	$7220 \pm 0$
	FF	1e-07		$.505 \pm .001$		$.364 \pm .035$	$.1 \pm .0$	$6195 \pm 0$
WI	FF	1e-06		$.503 \pm .010$		$.337 \pm .007$	$1.3 \pm 1.0$	$5169 \pm 0$
v V	SSD	1e-06		$.555 \pm .005$		$.982 \pm .005$	$479.0 \pm 4.0$	$28196 \pm 4$
	SSD	1e-05		$.553 \pm .005$		$.982 \pm .005$	$476.8 \pm 4.3$	$28728 \pm 6$
	SSD	1e-04	$.633 \pm .015$	$.552 \pm .004$	$.933 \pm .003$	$.982 \pm .005$	$477.4 \pm 4.5$	$28462 \pm 5$

Table 16: Comparison of Unlearners on BERT trained on IMDB (mean  $\pm$  std).

Group	Method	LR	LUMA			F1 (forget)		GPU (MB)
_	Orig.	-				$.997 \pm .002$	$4914.0 \pm 2965.0$	
_	Gold	-				$.928 \pm .008$	$4914.0 \pm 2965.0$	
	GD					$.999 \pm .001$	$1329.0 \pm 2.9$	$46742 \pm 6815$
	GD	1e-05	$.667 \pm .045$	$0.539 \pm 0.001$	$0.941 \pm 0.002$	$.996 \pm .004$	$1329.0 \pm 4.1$	$52408 \pm 4716$
	SRL	1e-06	$.653 \pm .047$	$.545 \pm .002$	$0.946 \pm 0.000$	$.999 \pm .001$	$1401.0 \pm .9$	$53724 \pm 4712$
	SRL		$.661 \pm .043$		1	$.994 \pm .002$	$1402.0 \pm .2$	$55041 \pm 4713$
FT	NG	1e-06	$.051 \pm .031$	$.497 \pm .000$	$0.408 \pm 0.090$	$.421 \pm .096$	$75.9 \pm 4.8$	$48925 \pm 6893$
ГІ	NG	1e-05	$.025 \pm .002$	$.501 \pm .004$	$0.335 \pm 0.000$	$.340 \pm .009$	$75.8 \pm 5.0$	$57670 \pm 4710$
	ANG	1e-06	$.660 \pm .056$	$.502 \pm .010$	$0.935 \pm 0.003$	$.953 \pm .014$	$271.0 \pm 8.4$	$58977 \pm 4704$
	ANG	1e-05	$.095 \pm .101$	$.527 \pm .040$	$.514 \pm .257$	$.415 \pm .134$	$271.0 \pm 8.5$	$60357 \pm 4709$
	UNSIR	1e-06	$\textbf{.649} \pm .049$	$.545 \pm .000$	$0.945 \pm 0.000$	$\textbf{.998} \pm .002$	$1478.0 \pm 5.9$	$61675 \pm 4702$
	UNSIR	1e-05	$\textbf{.656} \pm .036$	$.538 \pm .004$	$0.933 \pm 0.010$	$\textbf{.991} \pm .002$	$1478.0 \pm 6.1$	$62993 \pm 4700$
	CF-k	1e-06	$.716 \pm .038$	$.549 \pm .006$	$.938 \pm .005$	$.997 \pm .002$	$512.7 \pm 2.9$	$56362 \pm 4719$
	$\mathtt{CF} - k$	1e-05	$.715 \pm .039$	$.551 \pm .007$	$0.939 \pm 0.005$	$\textbf{.997} \pm .001$	$512.4 \pm 2.7$	$49558 \pm 5998$
	$\mathtt{EU-}k$	1e-06	$\textbf{.671} \pm .037$	$.554 \pm .002$	$0.937 \pm 0.005$	$\textbf{.997} \pm .001$	$1024.0 \pm 3.6$	$49558 \pm 5998$
SWM	$\mathtt{EU}{-}k$	1e-05	$\textbf{.}674 \pm .042$	$.551 \pm .002$	$.939 \pm .004$	$\textbf{.997} \pm .001$	$1025.0 \pm 2.5$	$48925 \pm 6893$
	SalUn	1e-06	$.580 \pm .000$	$.556 \pm .000$	$0.944 \pm 0.000$	$1.000\pm.000$	$1983.0 \pm .0$	$49122 \pm 0$
	SalUn	1e-05	$\textbf{.596} \pm .000$	$0.540 \pm 0.000$	$0.926 \pm 0.000$	$\textbf{.992} \pm .000$	$1988.0 \pm .0$	$61327 \pm 0$
	SalUn	1e-04	$.026 \pm .000$	$.498 \pm .000$	$0.335 \pm 0.000$	$\textbf{.344} \pm .000$	$1988.0 \pm .0$	$59980 \pm 0$
	BT	1e-06	$.656 \pm .000$	$.540 \pm .000$	$.890 \pm .000$	$.948 \pm .000$	$2427.0 \pm .0$	$67629 \pm 0$
	BT	1e-05	$.309 \pm .335$	$.535 \pm .017$	$.672 \pm .305$	$.748 \pm .340$	$2553.0 \pm 177.3$	$53771 \pm 21984$
DIS	SCRUB	1e-06	$\textbf{.628} \pm .064$	$.548 \pm .007$	$0.944 \pm 0.000$	$1.000\pm.000$	$1904.0 \pm 133.6$	$58627 \pm 17586$
	SCRUB	1e-05	$\textbf{.638} \pm .055$	$.543 \pm .001$	$.935 \pm .009$	$.993 \pm .007$	$1904.0 \pm 132.5$	$60541 \pm 16721$
	SCRUB	1e-04	$.024 \pm .000$	$.500 \pm .000$	$0.332 \pm 0.000$	$.323 \pm .000$	$1998.0 \pm .0$	$48219 \pm 0$
	FF	1e-08	$.065 \pm .047$	$.505 \pm .010$	$.447 \pm .097$	$.436 \pm .097$	$91.3 \pm 2.3$	$57216 \pm 16119$
	FF	1e-07	$.086 \pm .000$	$.498 \pm .002$	$0.500 \pm 0.011$	$.497 \pm .003$	$89.2 \pm 3.3$	$55539 \pm 16120$
WI	FF	1e-06	$.051 \pm .025$	$.498 \pm .003$	$1.423 \pm .063$	$.416 \pm .058$	$91.1 \pm 4.1$	$61832 \pm 16759$
W⊥	SSD	1e-06	$.642 \pm .065$	$.555 \pm .015$	$0.941 \pm 0.000$	$.999 \pm .001$	$1492.0 \pm 83.8$	$58894 \pm 16117$
	SSD	1e-05	$.641 \pm .063$	$0.556 \pm 0.013$	$0.941 \pm 0.000$	$.999 \pm .001$	$1492.0 \pm 83.1$	$61221 \pm 17715$
	SSD	1e-04	$\textbf{.642} \pm .062$	$0.555 \pm .012$	$0.941 \pm 0.000$	$\textbf{.999} \pm .001$	$1491.0 \pm 83.7$	$60790 \pm 17710$

Table 17: Comparison of Unlearners on  $MLP_{small}$  trained on Adult (mean  $\pm$  std).

							DunTime	
_	Method		LUMA			F1 (forget)	RunTime	GPU (MB)
_	Orig.	-				$.795 \pm .001$	$207.5 \pm 9.1$	$21 \pm 0$
_	Gold	-		$.499 \pm .001$			$207.5 \pm 9.1$	$21 \pm 0$
	GD	1e-04		$.499 \pm .002$		$.794 \pm .001$	$140.4 \pm 7.0$	22 ± 0
	GD	1e-03			$.793 \pm .001$		$138.5 \pm 2.5$	$22 \pm 0$
	GD	1e-02		$.499 \pm .000$		$.788 \pm .002$	$129.2 \pm 6.5$	$22 \pm 0$
	SRL	1e-04			$.788 \pm .002$		$163.8 \pm 11.9$	$23 \pm 0$
	SRL	1e-03		$.499 \pm .001$		$.795 \pm .001$	$179.4 \pm \textbf{2.6}$	$23 \pm \mathrm{0}$
	SRL	1e-02		$0.500 \pm 0.002$		$.790 \pm .002$	$166.1 \pm 3.7$	$23 \pm 0$
	NG	1e-04		$0.500 \pm 0.001$		$.732 \pm .015$	$23.4 \pm 1.1$	$25 \pm 0$
FΤ	NG		$1.107 \pm .072$			$.352 \pm .137$	$22.8 \pm .3$	$25 \pm 0$
	NG	1e-02		$.498 \pm .002$		$.432 \pm .000$	$25.5 \pm .1$	$25 \pm 0$
	ANG	1e-04	$0.580 \pm 0.010$	$.499 \pm .002$	$.788 \pm .002$	$.794 \pm .002$	$262.6 \pm 2.3$	$26 \pm 0$
	ANG	1e-03	$.129 \pm .012$	$.498 \pm .003$	$.419 \pm .015$	$.427 \pm .013$	$240.6 \pm \textbf{4.7}$	$26 \pm 0$
	ANG	1e-02	$.119 \pm .038$	$.498 \pm .001$	$.405 \pm .048$	$.411 \pm .047$	$268.6 \pm 1.3$	$26 \pm 0$
	UNSIR	1e-04	$.654 \pm .013$	$.498 \pm .002$	$.789 \pm .002$	$.791 \pm .002$	$174.7 \pm 5.8$	$27 \pm 0$
	UNSIR	1e-03	$.646 \pm .006$	$.498 \pm .001$	$.793 \pm .001$	$.796 \pm .001$	$184.8 \pm \textbf{4.7}$	$27\pm o$
	UNSIR	1e-02	$.657 \pm .005$	$.499 \pm .001$	$.790 \pm .001$	$.788 \pm .001$	$171.9 \pm \text{1.5}$	$26 \pm 0$
	CF-k	1e-04	$.710 \pm .001$	$.499 \pm .001$	$.792 \pm .001$	$.792 \pm .001$	$130.2 \pm 6.0$	24 ± 0
	CF-k	1e-03	$.695 \pm .002$	$0.500 \pm 0.000$	$.794 \pm .000$	$.795 \pm .001$	$142.6 \pm 7.2$	$24 \pm 0$
	CF-k	1e-02	$.705 \pm .008$	$.498 \pm .000$	$.793 \pm .001$	$.797 \pm .000$	$134.0 \pm 2.3$	$24 \pm 0$
	EU- $k$	1e-04	$.270 \pm .009$	$.498 \pm .001$	$.788 \pm .002$	$.789 \pm .003$	$1333.0 \pm \textbf{14.0}$	$24 \pm 0$
SWM	EU- $k$	1e-03	$.268 \pm .005$	$.499 \pm .001$	$.793 \pm .002$	$.797 \pm .002$	$1359.0 \pm 4.1$	$24 \pm 0$
	EU- $k$	1e-02	$.267 \pm .005$	$.498 \pm .001$	$.791 \pm .002$	$.798 \pm .002$	$1367.0 \pm \textbf{39.5}$	$24 \pm 0$
	SalUn	1e-04	$.689 \pm .005$	$0.500 \pm 0.001$	$.788 \pm .001$	$.789 \pm .000$	$137.7 \pm 6.9$	$32 \pm 0$
	SalUn	1e-03	$.691 \pm .015$	$0.500 \pm 0.002$	$.790 \pm .002$	$.794 \pm .001$	$139.7 \pm \text{14.6}$	$31 \pm 0$
	SalUn	1e-02	$.698 \pm .008$	$.499 \pm .001$	$.789 \pm .002$	$.789 \pm .001$	$131.6 \pm 9.2$	$31 \pm 0$
-	BT	1e-04	$.339 \pm .078$	$.501 \pm .003$	$.593 \pm .053$	$.595 \pm .051$	$224.5 \pm 2.9$	$28 \pm 0$
	BT	1e-03	$.476 \pm .066$	$0.500 \pm 0.001$	$.686 \pm .045$	$.692 \pm .044$	$234.9 \pm \textbf{3.3}$	$28 \pm 0$
DIG	BT	1e-02	$.256 \pm .120$	$.498 \pm .001$	$.527 \pm .090$	$.534 \pm .093$	$224.9 \pm 7.9$	$27\pm 0$
DIS	SCRUB	1e-04	$.716 \pm .004$	$.499 \pm .002$	$.792 \pm .001$	$.793 \pm .001$	$120.9 \pm 3.1$	$29 \pm 0$
	SCRUB	1e-03	$.715 \pm .008$	$0.500 \pm 0.000$	$.793 \pm .002$	$.796 \pm .001$	$121.5 \pm 4.2$	$29 \pm 0$
	SCRUB	1e-02	$.692 \pm .010$	$.498 \pm .002$	$.787 \pm .006$	$.786 \pm .003$	$134.8 \pm .4$	$29 \pm \mathrm{o}$
	FF	1e-08	$.949 \pm .009$	$.499 \pm .001$	$.784 \pm .001$	$.791 \pm .002$	$5.1 \pm .7$	18±0
	FF	1e-07	$.941 \pm .017$	$.501 \pm .001$	$.785 \pm .002$	$.785 \pm .005$	$5.6 \pm 1.5$	$18 \pm 0$
T-7 T	FF	1e-06	$.941 \pm .020$	$1.500 \pm .001$	$1.782 \pm .009$	$.786 \pm .010$	$5.1 \pm 1.5$	$18 \pm 0$
WI	SSD	1e-04			$.794 \pm .002$	$.795 \pm .001$	$139.3 \pm 4.0$	$31 \pm 0$
	SSD	1e-03		$.500 \pm .001$		$.795 \pm .001$	$137.3 \pm 4.6$	$30 \pm 0$
	SSD	1e-02		$.499 \pm .000$		$.795 \pm .001$	$129.9 \pm 1.8$	$30 \pm 0$
						•		

Table 18: Comparison of Unlearners on  $MLP_{large}$  trained on Adult (mean  $\pm$  std).

Group	Method	LR	LUMA			F1 (forget)		GPU (MB
_	Orig.	-	$.631 \pm .000$	$\textbf{.}498 \pm .001$	$.793 \pm .002$	$.798 \pm .001$	$135.7 \pm 7.9$	$32 \pm 0$
_	Gold	-	$.636 \pm .000$	$.499 \pm .001$	$.791 \pm .002$	$.792 \pm .001$	$135.7 \pm 7.9$	$32 \pm 0$
	GD	1e-04	$.707 \pm .003$	$.499 \pm .001$	$.795 \pm .001$	$.798 \pm .000$	$87.6 \pm 3.6$	36±0
	GD	1e-03	$.710 \pm .001$	$\textbf{.}498 \pm .002$	$.791 \pm .003$	$.795 \pm .002$	$88.3 \pm 4.9$	$35\pm 0$
	GD	1e-02	$.709 \pm .003$	$.500 \pm .001$	$.786 \pm .004$	$.788 \pm .002$	$87.2 \pm 5.3$	$33\pm 0$
	SRL	1e-04	$.667 \pm .003$	$.498 \pm .001$	$.788 \pm .005$	$.791 \pm .000$	$11.0 \pm 5.8$	$40 \pm {\scriptscriptstyle 0}$
	SRL	1e-03	$.668 \pm .002$	$.500 \pm .000$	$.790 \pm .002$	$.794 \pm .002$	$11.2 \pm 4.7$	$39\pm 0$
	SRL	1e-02	$.667 \pm .002$	$\textbf{.}498 \pm .001$	$.789 \pm .001$	$.789 \pm .002$	$11.6 \pm 5.6$	$37\pm 0$
	NG	1e-04	$.096 \pm .085$	$.499 \pm .002$	$.330 \pm .136$	$.335 \pm .135$	$16.0 \pm .2$	$47\pm 0$
FT	NG	1e-03	$.150 \pm .001$	$.500 \pm .002$	$.434 \pm .000$	$.432 \pm .000$	$16.2 \pm .4$	$46\pm 0$
	NG	1e-02	$.108 \pm .073$	$\textbf{.}499 \pm .001$	$.352 \pm .141$	$.352 \pm .137$	$16.4 \pm .3$	$44\pm0$
	ANG	1e-04	$.591 \pm .006$	$\textbf{.}499 \pm .000$	$.767 \pm .004$	$.773 \pm .003$	$140.7 \pm 1.5$	$50 \pm 0$
	ANG	1e-03	$.117 \pm .014$	$\textbf{.}499 \pm .001$	$.403 \pm .018$	$.409 \pm .017$	$144.6 \pm 2.3$	$49\pm 0$
	ANG	1e-02	$.062 \pm .066$	$.501 \pm .004$	$.271 \pm .141$	$.273 \pm .137$	$148.6 \pm 1.1$	$48 \pm 0$
	UNSIR	1e-04	$.666 \pm .010$	$\textbf{.}499 \pm .001$	$.791 \pm .002$	$.797 \pm .001$	$105.6 \pm .0$	$54 \pm 0$
	UNSIR	1e-03	$.668 \pm .008$	$\textbf{.}498 \pm .001$	$.792 \pm .003$	$.795 \pm .002$	$105.2 \pm .6$	$53\pm 0$
	UNSIR	1e-02	$.668 \pm .008$	$\textbf{.}499 \pm .001$	$.787 \pm .004$	$.790 \pm .002$	$105.2 \pm .8$	$51\pm 0$
	CF-k	1e-04	$.711 \pm .007$	$.499 \pm .000$	$.791 \pm .003$	$.797 \pm .001$	$85.5 \pm 1.1$	41 ± 0
	CF-k	1e-03	$.710 \pm .007$	$\textbf{.}500 \pm .001$	$.793 \pm .003$	$.797 \pm .002$	$85.2 \pm .9$	$41\pm {\rm 0}$
	CF-k	1e-02	$.706 \pm .002$	$\textbf{.}498 \pm .001$	$.794 \pm .003$	$.797 \pm .000$	$87.7 \pm 5.3$	$41\pm {\rm 0}$
	EU-k	1e-04	$.240 \pm .012$	$\textbf{.}499 \pm .001$	$.789 \pm .003$	$.793 \pm .004$	$850.2 \pm 6.4$	$44 \pm 0$
SWM	EU-k	1e-03	$.242 \pm .010$	$\textbf{.}498 \pm .001$	$.793 \pm .002$	$.796 \pm .002$	$848.4 \pm 15.6$	
	EU-k	1e-02	$.241 \pm .009$	$\textbf{.500} \pm .000$	$.794 \pm .004$	$.796 \pm .001$	$856.0 \pm 16.2$	$42 \pm 0$
	SalUn	1e-04	$.676 \pm .014$	$\textbf{.}500 \pm .001$	$.786 \pm .003$	$.790 \pm .001$	$100.8 \pm 11.9$	$57 \pm 27$
	SalUn	1e-03	$.675 \pm .018$	$\textbf{.}499 \pm .000$	$.791 \pm .001$	$.796 \pm .002$	$103.1 \pm 14.3$	$55 \pm 27$
	SalUn	1e-02	$.676 \pm .016$	$\textbf{.}498 \pm .001$	$.792 \pm .001$	$.791 \pm .003$	$103.7 \pm 14.3$	$53 \pm 26$
	BT	1e-04	$.354 \pm .260$	$.499 \pm .003$	$.576 \pm .247$	$.578 \pm .238$	$155.1 \pm 1.6$	$58 \pm 0$
	BT	1e-03	$.432 \pm .154$	$\textbf{.}504 \pm .005$	$.660 \pm .106$	$.669 \pm .112$	$156.1 \pm .9$	$57 \pm \mathrm{0}$
DIC	BT	1e-02	$.427 \pm .175$	$.502 \pm .004$	$.661 \pm .127$	$.663 \pm .134$	$155.5 \pm 1.8$	$56 \pm 0$
DIS	SCRUB	1e-04	$.708 \pm .009$	$\textbf{.}498 \pm .001$	$.795 \pm .001$	$.799 \pm .000$	$80.8 \pm .2$	$61\pm 0$
	SCRUB	1e-03	$.711 \pm .007$	$\textbf{.}498 \pm .002$	$.789 \pm .003$	$.791 \pm .004$	$81.4 \pm 1.0$	$60 \pm 0$
	SCRUB	1e-02	$.708 \pm .012$	$\textbf{.}499 \pm .003$	$.790 \pm .005$	$.791 \pm .007$	$82.1 \pm .4$	$60 \pm 0$
	FF	1e-08	$.897 \pm .046$	$.501 \pm .002$	$.773 \pm .018$	$.770 \pm .021$	$8.3 \pm .8$	$21 \pm 0$
	FF	1e-07	$.936 \pm .013$	$\textbf{.}500 \pm .001$	$.786 \pm .007$	$.785 \pm .006$	$8.0 \pm .8$	$20\pm 0$
WI	FF	1e-06	$.897 \pm .048$	$\textbf{.}500 \pm .002$	$.768 \pm .017$	$.770 \pm .020$	$7.2 \pm .4$	$19\pm 0$
VV T	SSD		$.676 \pm .023$			$.798 \pm .001$	$103.7 \pm 19.6$	$52 \pm {\scriptstyle 26}$
	SSD	1e-03	$.675 \pm .023$	$\textbf{.499} \pm .001$	$.793 \pm .002$	$.798 \pm .001$	$104.4 \pm 2.0$	$51\pm 26$
	SSD	1e-02	$.697 \pm .016$	$\textbf{.}499 \pm .001$	$.793 \pm .002$	$.798 \pm .001$	$91.5 \pm .9$	$50 \pm 26$

Table 19: Comparison of Unlearners on  $MLP_{small}$  trained on Spotify (mean  $\pm$  std).

				earners on M			-	
Group	Method	LR	LUMA	UMIA	F1 (test)	F1 (forget)		
_	Orig.	-		$0.499 \pm 0.004$				$19 \pm 0$
_	Gold	-		$0.499 \pm 0.006$				19±0
	GD	l		$0.500 \pm 0.004$				19±0
	GD	l		$0.498 \pm 0.003$				$19 \pm 0$
	GD			$0.497 \pm 0.002$				$18 \pm 0$
	SRL	l		$0.497 \pm 0.004$			1	$20 \pm 0$
	SRL	l		$0.499 \pm 0.001$				$19 \pm 0$
	SRL	l		$0.496 \pm 0.001$			$6.6 \pm 0.3$	$18 \pm 0$
	NG			$0.498 \pm 0.002$			$2.4 \pm 0.4$	$20 \pm 0$
FT	NG	1e-03	$0.836 \pm 0.038$	$0.496 \pm 0.006$	$0.551 \pm 0.018$	$0.546 \pm 0.006$		$20 \pm 0$
	NG			$0.501 \pm 0.011$				$18 \pm 0$
	ANG			$0.497 \pm 0.001$				$21 \pm 0$
	ANG	l		$0.502 \pm 0.011$			1	$21 \pm 0$
	ANG			$0.529 \pm 0.007$			$10.2 \pm 0.3$	$18 \pm 0$
	UNSIR			$0.499 \pm 0.004$			$15.7 \pm 0.8$	$21 \pm 0$
	UNSIR	1e-03	$0.884 \pm 0.005$	$0.497 \pm 0.002$	$0.609 \pm 0.005$	$0.615 \pm \textbf{0.005}$		$21 \pm 0$
	UNSIR			$0.495 \pm 0.002$				$19 \pm 0$
	CF-k	l		$0.498 \pm 0.001$				$20 \pm 0$
	CF-k			$0.499 \pm 0.000$				$20 \pm 0$
	CF-k			$0.501 \pm 0.002$			$5.2 \pm 0.3$	$18 \pm 0$
	EU- $k$			$0.496 \pm 0.004$				$20 \pm 0$
SWM	EU- $k$			$0.497 \pm 0.001$				$20 \pm 0$
	EU- $k$			$0.497 \pm 0.003$				$18 \pm 0$
	SalUn			$0.501 \pm 0.002$				$23 \pm 0$
	SalUn			$0.498 \pm 0.003$				$23 \pm 0$
	SalUn			$0.496 \pm 0.001$				$19 \pm 0$
	BT			$0.501 \pm 0.005$				$22 \pm 0$
	BT			$0.506 \pm 0.007$				$21 \pm 0$
DIS	BT			$0.505 \pm 0.002$				$19 \pm 0$
DIO	SCRUB			$0.496 \pm 0.002$				$22 \pm 0$
	SCRUB			$0.498 \pm 0.001$				$22 \pm 0$
	SCRUB			$0.499 \pm 0.001$				$19 \pm 0$
	FF	l		$0.497 \pm 0.001$			$15.1 \pm 0.8$	$18 \pm 0$
	FF			$0.497 \pm 0.006$			$15.5 \pm 0.2$	$18 \pm 0$
WI	FF	l		$0.497 \pm 0.006$			$15.8 \pm 0.1$	$18 \pm 0$
V V	SSD			$0.501 \pm 0.004$				$22 \pm 0$
	SSD			$0.500 \pm 0.006$				$22 \pm 0$
	SSD	1e-02	$0.901 \pm 0.004$	$0.501 \pm 0.003$	$0.612 \pm 0.002$	$0.620 \pm \textbf{0.004}$	$8.0 \pm 0.6$	$19 \pm 0$
				l	l			

Table 20: Comparison of Unlearners on  $MLP_{large}$  trained on Spotify (mean  $\pm$  std).

			arison of Uni					
Group	Method	LR	LUMA	UMIA	F1 (test)	F1 (forget)		
_	Orig.	-			$0.635 \pm 0.009$			$30 \pm 0$
_	Gold				$0.629 \pm 0.002$		$119.6 \pm 1.7$	30 ± 0
	GD		$0.810 \pm 0.003$				$6.0 \pm 0.1$	33 ± 0
	GD		$0.850 \pm 0.022$				$6.0 \pm 0.1$	$31 \pm 0$
	GD		$0.915 \pm 0.018$				$4.7 \pm 0.1$	$22 \pm 0$
	SRL		$0.796 \pm 0.006$				$7.6 \pm 0.1$	$36 \pm 0$
	SRL		$0.842 \pm 0.012$				$7.8 \pm 0.1$	$35 \pm 0$
	SRL		$0.764 \pm 0.002$				$6.8 \pm 0.0$	$23 \pm 0$
FT	NG		$0.875 \pm 0.006$				$1.5 \pm 0.0$	$42 \pm 0$
	NG		$0.023 \pm 0.000$	l			$1.5 \pm 0.0$	$40 \pm 0$
	NG		$0.023 \pm 0.000$				$1.3 \pm 0.0$	$25 \pm 0$
	ANG		$0.875 \pm 0.002$				$11.9 \pm 0.1$	$45 \pm 0$
	ANG		$0.085 \pm 0.021$					$44 \pm 0$
	ANG		$0.024 \pm 0.001$		$0.024 \pm 0.002$		$10.7 \pm 0.3$	$26 \pm 0$
	UNSIR		$0.792 \pm 0.004$				$9.3 \pm 0.0$	$48 \pm 0$
	UNSIR		$0.824 \pm 0.011$	l	$0.642 \pm 0.008$		$9.4 \pm 0.1$	$47 \pm 0$
	UNSIR		$0.885 \pm 0.012$				$8.5 \pm 0.2$	$27 \pm 0$
	CF-k		$0.810 \pm 0.010$				$5.9 \pm 0.1$	$37 \pm 0$
	CF-k		$0.809 \pm 0.011$				$6.0 \pm 0.1$	$38 \pm 0$
	CF-k				$0.642 \pm 0.007$		$5.3 \pm 0.1$	$23 \pm 0$
	$_{ m EU-}k$		$0.687 \pm 0.009$				$59.6 \pm 0.5$	$39 \pm 0$
SWM	EU- $k$		$0.692 \pm 0.000$				$59.6 \pm 0.3$	$39 \pm 0$
	EU- $k$		$0.687 \pm 0.009$				$52.8 \pm 0.2$	$24 \pm 0$
	SalUn		$0.786 \pm 0.009$				$8.0 \pm 0.6$	$63 \pm 0$
	SalUn		$0.817 \pm 0.021$				$8.3 \pm 0.6$	$60 \pm 0$
	SalUn	1e-02	$0.779 \pm 0.026$	$0.501 \pm 0.004$	$0.564 \pm 0.006$	$0.562 \pm 0.014$	$8.4 \pm 0.3$	$31 \pm 0$
	BT		$0.785 \pm 0.029$		$0.583 \pm 0.027$		$13.2 \pm 1.0$	$51 \pm 0$
	BT		$0.641 \pm 0.114$				$13.6 \pm 0.8$	$50 \pm 0$
DIS	BT		$0.667 \pm 0.060$				$12.8 \pm 0.7$	$28 \pm 0$
DID	SCRUB		$0.796 \pm 0.013$				$7.6 \pm 0.1$	$54 \pm 0$
	SCRUB		$0.841 \pm 0.004$				$7.6 \pm 0.1$	$54 \pm 0$
	SCRUB	1e-02	$0.325 \pm 0.072$	$0.498 \pm 0.000$	$0.382 \pm 0.035$	$0.377 \pm 0.041$	$6.8 \pm 0.1$	$29 \pm 0$
	FF		$0.808 \pm 0.020$				$19.3 \pm 0.6$	$21 \pm 0$
	FF		$0.812 \pm 0.017$				$19.4 \pm 0.2$	$20 \pm 0$
WI	FF	1e-06	$0.769 \pm 0.107$	$0.513 \pm 0.009$	$0.572 \pm 0.051$	$0.605 \pm \textbf{0.063}$	$19.3 \pm 0.2$	$19 \pm 0$
AA T	SSD		$0.807 \pm 0.018$				$8.7 \pm 0.2$	$59 \pm 0$
	SSD		$0.806 \pm 0.013$				$8.9 \pm 0.1$	$58 \pm 0$
	SSD	1e-02	$0.827 \pm 0.017$	$0.518 \pm 0.007$	$0.635 \pm 0.009$	$0.689 \pm \textbf{0.008}$	$8.1 \pm 0.1$	$30 \pm 0$
					l			

Table 21: Comparison of Unlearners on  $GCN_{small}$  trained on BACE (mean  $\pm$  std).

					$CN_{small}$ tra			
Group	Method		LUMA	UMIA	F1 (test)			GPU (MB)
_	Orig.	-			$0.572 \pm 0.009$			18±0
_	Gold	-			$0.540 \pm 0.019$			18±0
	GD				$0.570 \pm 0.007$			18 ± 0
	GD				$0.566 \pm 0.009$			$18 \pm 0$
	GD				$0.552 \pm 0.024$			$18 \pm 0$
	SRL				$0.572 \pm 0.004$			$18 \pm 0$
	SRL	1			$0.536 \pm 0.005$		I .	$18 \pm 0$
	SRL	1			$0.355 \pm 0.000$		I .	$18 \pm 0$
	NG				$0.414 \pm 0.004$			$18 \pm 0$
FT	NG				$0.355 \pm 0.000$			$18 \pm 0$
	NG				$0.355 \pm 0.000$			$18 \pm 0$
	ANG				$0.537 \pm 0.008$			$18 \pm 0$
	ANG	1			$0.379 \pm 0.008$		I .	$18 \pm 0$
	ANG	1			$0.355 \pm 0.000$			$18 \pm 0$
	UNSIR				$0.570 \pm 0.007$			$18 \pm 0$
	UNSIR				$0.566 \pm 0.009$		I .	$18 \pm 0$
	UNSIR				$0.552 \pm 0.024$			$18 \pm 0$
	CF-k				$0.570 \pm 0.007$			$18 \pm 0$
	CF-k				$0.566 \pm 0.009$			$18 \pm 0$
	CF-k				$0.552 \pm 0.024$			$18 \pm 0$
	EU-k				$0.573 \pm 0.011$			$18 \pm 0$
SWM	EU-k				$0.577 \pm 0.007$			$18 \pm 0$
	EU-k				$0.624 \pm 0.018$			$18 \pm 0$
	SalUn				$0.572 \pm 0.004$			$19 \pm 0$
	SalUn				$0.539 \pm 0.005$			$19 \pm 0$
	SalUn				$0.355 \pm 0.000$			$19 \pm 0$
	BT				$0.588 \pm 0.008$			$19 \pm 0$
	BT				$0.589 \pm 0.022$		$3.2 \pm 0.0$	$19 \pm 0$
DIS	BT	1			$0.503 \pm 0.010$		$3.2 \pm 0.0$	$19 \pm 0$
DIO	SCRUB				$0.570 \pm 0.007$			$19 \pm 0$
	SCRUB				$0.566 \pm 0.009$			$19 \pm 0$
	SCRUB				$0.552 \pm 0.024$			$19 \pm 0$
	FF	1			$0.417 \pm 0.048$		I .	$19 \pm 0$
	FF				$0.498 \pm 0.117$			$19 \pm 0$
WI	FF				$0.523 \pm 0.125$			$19 \pm 0$
VV I	SSD				$0.572 \pm 0.009$			$19 \pm 0$
	SSD				$0.572 \pm 0.009$		$1.9 \pm 0.0$	$19 \pm 0$
	SSD	1e-02	$0.865 \pm 0.094$	$0.527 \pm 0.023$	$0.572 \pm 0.009$	$0.662 \pm \textbf{0.011}$	$1.9 \pm 0.0$	$19 \pm 0$
					l		·	

Table 22: Comparison of Unlearners on  $GCN_{large}$  trained on BACE (mean  $\pm$  std).

						Thea on BAC.		
Group	Method		LUMA	UMIA	F1 (test)	F1 (forget)		
_	Orig.	-				$0.711 \pm 0.002$		$18 \pm 0$
_	Gold	-				$0.679 \pm 0.026$		18 ± 0
	GD			$0.519 \pm 0.017$			$3.1 \pm 4.0$	19 ± 0
	GD			$0.534 \pm 0.008$			$3.0 \pm 4.0$	$19 \pm 0$
	GD					$0.468 \pm 0.180$		$18 \pm 0$
	SRL			$0.536 \pm 0.017$			$4.0 \pm 5.1$	$19 \pm 0$
	SRL			$0.504 \pm 0.031$			$4.1 \pm 5.2$	$19 \pm 0$
	SRL			$0.510 \pm 0.022$			$3.9 \pm 5.0$	$19 \pm 0$
	NG			$0.488 \pm 0.011$			$0.2 \pm 0.0$	$19 \pm 0$
FT	NG	1e-03	$0.283 \pm 0.060$	$0.499 \pm 0.008$	$0.355 \pm 0.000$	$0.348 \pm 0.000$	$0.2 \pm 0.0$	$19 \pm 0$
	NG			$0.506 \pm 0.009$			$0.2 \pm 0.0$	$19 \pm 0$
	ANG	1e-04	$0.833 \pm 0.099$	$0.505 \pm 0.020$	$0.620 \pm 0.009$	$0.702 \pm \textbf{0.003}$	$0.7 \pm 0.0$	$20 \pm 0$
	ANG	1e-03	$0.869 \pm 0.031$	$0.496 \pm 0.017$	$0.534 \pm 0.028$	$0.658 \pm \textbf{0.015}$	$0.7 \pm 0.0$	$19 \pm 0$
	ANG	1e-02	$0.284 \pm 0.063$	$0.503 \pm 0.007$	$0.355 \pm 0.000$	$0.348 \pm \textbf{0.000}$	$0.7 \pm 0.0$	$19 \pm 0$
	UNSIR			$0.513 \pm 0.029$			$1.0 \pm 0.0$	$20 \pm 0$
	UNSIR	1e-03	$0.850 \pm 0.097$	$0.528 \pm 0.006$	$0.617 \pm 0.009$	$0.707 \pm 0.005$	$1.0 \pm 0.0$	$20 \pm 0$
	UNSIR			$0.512 \pm 0.014$			$1.0 \pm 0.0$	$20 \pm 0$
	CF-k			$0.522 \pm 0.017$			$3.2 \pm 4.1$	19 ± 0
	$\mathtt{CF-}k$			$0.539 \pm 0.019$			$3.2 \pm 4.1$	$19 \pm 0$
	$\mathtt{CF-}k$					$0.468 \pm 0.180$		$19 \pm 0$
	EU- $k$					$0.696 \pm \textbf{0.002}$		$19 \pm 0$
SWM	EU- $k$					$0.710 \pm \textbf{0.005}$		$19 \pm 0$
	EU- $k$					$0.577 \pm \textbf{0.157}$		$19 \pm 0$
	SalUn					$0.710 \pm \textbf{0.002}$		$20 \pm 0$
	SalUn			$0.536 \pm 0.023$			$1.3 \pm 0.0$	$21 \pm 0$
	SalUn			$0.521 \pm 0.036$			$1.3 \pm 0.0$	$20 \pm 0$
	BT			$0.531 \pm 0.009$			$2.1 \pm 0.0$	$20 \pm 0$
	BT			$0.529 \pm 0.041$			$2.2 \pm 0.1$	$20 \pm 0$
DIS	BT	1e-02	$0.803 \pm 0.062$	$0.515 \pm 0.019$	$0.550 \pm 0.019$	$0.604 \pm \textbf{0.034}$	$2.1 \pm 0.0$	$20 \pm 0$
טוט	SCRUB			$0.523 \pm 0.013$			$1.4 \pm 0.0$	$20 \pm 0$
	SCRUB			$0.506 \pm 0.011$			$1.5 \pm 0.0$	$20 \pm 0$
	SCRUB			$0.495 \pm 0.006$			$1.4 \pm 0.0$	$20 \pm 0$
	FF			$0.508 \pm 0.013$			$6.8 \pm 0.1$	$20 \pm 0$
	FF			$0.514 \pm 0.008$			$6.7 \pm 0.0$	$20 \pm 0$
WI	FF	1e-06	$0.830 \pm 0.108$	$0.526 \pm 0.017$	$0.607 \pm 0.020$	$0.683 \pm \textbf{0.013}$	$6.9 \pm 0.1$	$20\pm 0$
VV I	SSD			$0.527 \pm 0.009$			$1.2 \pm 0.0$	$20\pm 0$
	SSD			$0.519 \pm 0.024$			$1.2 \pm 0.0$	$20\pm \mathrm{o}$
	SSD	1e-02	$0.803 \pm 0.084$	$0.522 \pm 0.031$	$0.630 \pm 0.004$	$0.711 \pm \textbf{0.002}$	$1.2 \pm 0.0$	$20\pm 0$
					<u> </u>			

Table 23: Comparison of Unlearners on  $GCN_{small}$  trained on BBBP (mean  $\pm$  std).

					$CN_{small}$ tra			
Group	Method		LUMA	UMIA	F1 (test)			GPU (MB)
_	Orig.	-			$0.665 \pm 0.014$			18±0
_	Gold	-			$0.681 \pm 0.012$			18±0
	GD				$0.671 \pm 0.010$			18 ± 0
	GD	l			$0.682 \pm 0.009$		$1.7 \pm 0.0$	$18 \pm 0$
	GD				$0.674 \pm 0.024$		$1.7 \pm 0.0$	$18 \pm 0$
	SRL				$0.684 \pm 0.010$		$2.4 \pm 0.0$	$18 \pm 0$
	SRL	l			$0.671 \pm 0.015$		$2.4 \pm 0.1$	$18 \pm 0$
	SRL	l			$0.538 \pm 0.091$		$2.2 \pm 0.0$	$18 \pm 0$
	NG				$0.566 \pm 0.083$		$0.5 \pm 0.0$	$18 \pm 0$
FT	NG	l			$0.435 \pm 0.000$		$0.6 \pm 0.2$	$18 \pm 0$
	NG				$0.435 \pm 0.000$			$18 \pm 0$
	ANG	l			$0.652 \pm 0.028$		$1.3 \pm 0.1$	$18 \pm 0$
	ANG				$0.484 \pm 0.045$		$1.5 \pm 0.2$	$18 \pm 0$
	ANG	l			$0.435 \pm 0.000$		$1.4 \pm 0.2$	$18 \pm 0$
	UNSIR	l			$0.671 \pm 0.010$		$2.1 \pm 0.0$	$18 \pm 0$
	UNSIR	l		l	$0.682 \pm 0.009$		$2.1 \pm 0.0$	$18 \pm 0$
	UNSIR				$0.674 \pm 0.024$		$2.1 \pm 0.0$	$18 \pm 0$
	CF-k				$0.671 \pm 0.010$			18 ± 0
	CF-k				$0.682 \pm 0.009$		$1.9 \pm 0.0$	$18 \pm 0$
	CF-k				$0.674 \pm 0.024$		$1.9 \pm 0.0$	$18 \pm 0$
	$_{ m EU-}k$	l			$0.676 \pm 0.016$			$18 \pm 0$
SWM	EU-k				$0.698 \pm 0.004$			$18 \pm 0$
	EU- $k$				$0.692 \pm 0.003$			$18 \pm 0$
	SalUn				$0.684 \pm 0.010$			$19 \pm 0$
	SalUn				$0.670 \pm 0.014$		$2.7 \pm 0.1$	$19 \pm 0$
	SalUn				$0.535 \pm 0.088$			$19 \pm 0$
	BT	l			$0.615 \pm 0.005$		$4.2 \pm 0.1$	$19 \pm 0$
	BT				$0.612 \pm 0.091$		$4.2 \pm 0.0$	$18 \pm 0$
DIS	BT	l	$0.650 \pm 0.297$		$0.554 \pm 0.114$		$4.3 \pm 0.0$	$18 \pm 0$
DIO	SCRUB				$0.671 \pm 0.010$			$19 \pm 0$
	SCRUB				$0.682 \pm 0.009$		$3.1 \pm 0.0$	$19 \pm 0$
	SCRUB				$0.674 \pm 0.024$		$3.1 \pm 0.0$	$19 \pm 0$
	FF				$0.681 \pm 0.015$		$10.3 \pm 0.1$	$19 \pm 0$
	FF				$0.576 \pm 0.122$		$10.2 \pm 0.1$	$19 \pm 0$
WI	FF	1e-06	$0.882 \pm 0.027$	$0.493 \pm 0.007$	$0.661 \pm 0.028$	$0.661 \pm 0.016$		$19 \pm 0$
VV I	SSD				$0.665 \pm 0.014$			$19 \pm 0$
	SSD				$0.665 \pm 0.014$			$19 \pm 0$
	SSD	1e-02	$0.927 \pm 0.031$	$0.496 \pm 0.006$	$0.665 \pm 0.014$	$0.676 \pm \textbf{0.013}$	$2.6 \pm 0.0$	$19 \pm 0$
					l			

Table 24: Comparison of Unlearners on  $GCN_{large}$  trained on BBBP (mean  $\pm$  std).

Group   I	· / · · · · · · · · · · · · · · · · · ·							~
_		LR	LUMA	UMIA	F1 (test)	F1 (forget)		, ,
_ (	Orig.	-				$0.705 \pm 0.007$		$18 \pm 0$
_	Gold	-			$0.712 \pm 0.007$		$55.1 \pm 0.2$	18 ± 0
	GD		$0.924 \pm 0.014$				$1.0 \pm 0.0$	19 ± 0
	GD		$0.921 \pm 0.034$				$1.0 \pm 0.0$	$19 \pm 0$
	GD					$0.657 \pm \textbf{0.030}$	$1.0 \pm 0.0$	$18 \pm 0$
	SRL		$0.930 \pm 0.013$				$1.4 \pm 0.0$	$19 \pm 0$
	SRL		$0.733 \pm 0.309$				$1.5 \pm 0.0$	$19 \pm 0$
	SRL		$0.777 \pm 0.127$				$1.3 \pm 0.0$	$19 \pm 0$
	NG		$0.408 \pm 0.126$				$0.3 \pm 0.0$	$19 \pm 0$
FT	NG	1e-03	$0.285 \pm 0.007$	$0.498 \pm 0.007$	$0.435 \pm 0.000$	$0.432 \pm 0.000$	$0.3 \pm 0.0$	$19 \pm 0$
	NG		$0.284 \pm 0.007$				$0.3 \pm 0.0$	$19 \pm 0$
	ANG		$0.916 \pm 0.022$				$0.9 \pm 0.0$	$19 \pm 0$
	ANG		$0.311 \pm 0.031$				$0.9 \pm 0.0$	$19 \pm 0$
	ANG		$0.284 \pm 0.007$				$0.9 \pm 0.0$	$19 \pm 0$
	UNSIR		$0.928 \pm 0.009$				$1.3 \pm 0.0$	$20 \pm 0$
I .	UNSIR		$0.917 \pm 0.039$				$1.4 \pm 0.0$	$20 \pm 0$
1	UNSIR		$0.866 \pm 0.111$				$1.3 \pm 0.0$	19 ± 0
	CF-k		$0.920 \pm 0.008$				$1.1 \pm 0.0$	19 ± 0
	CF-k		$0.917 \pm 0.034$				$1.1 \pm 0.0$	$19 \pm 0$
	CF-k		$0.863 \pm 0.110$				$1.2 \pm 0.0$	$19 \pm 0$
	EU- $k$		$0.713 \pm 0.005$				$34.0 \pm 0.0$	$19 \pm 0$
SWM	EU- $k$		$0.699 \pm 0.021$				$34.1 \pm 0.1$	$19 \pm 0$
	EU- $k$		$0.704 \pm 0.011$				$34.3 \pm 0.1$	$19 \pm 0$
I .	SalUn		$0.919 \pm 0.012$				$1.7 \pm 0.1$	$21 \pm 0$
I .	SalUn		$0.740 \pm 0.328$				$1.6 \pm 0.1$	$20 \pm 0$
	SalUn		$0.689 \pm 0.131$				$1.7 \pm 0.0$	18 ± 0
	BT		$0.879 \pm 0.028$				$2.9 \pm 0.0$	$20 \pm 0$
	BT		$0.854 \pm 0.043$				$2.8 \pm 0.0$	$20 \pm 0$
DIS	BT		$0.438 \pm 0.131$				$3.0 \pm 0.0$	$20 \pm 0$
	SCRUB		$0.917 \pm 0.010$				$1.9 \pm 0.0$	$20 \pm 0$
	SCRUB		$0.920 \pm 0.032$				$1.9 \pm 0.0$	$20\pm0$
	SCRUB		$0.864 \pm 0.113$				$1.9 \pm 0.0$	20 ± 0
	FF		$0.771 \pm 0.167$				$7.1 \pm 0.1$	$20 \pm 0$
	FF		$0.787 \pm 0.075$				$7.1 \pm 0.0$	$20 \pm 0$
WI	FF		$0.793 \pm 0.147$				$7.0 \pm 0.0$	$20 \pm 0$
** +	SSD		$0.926 \pm 0.004$				$1.5 \pm 0.1$	$20 \pm 0$
	SSD		$0.930 \pm 0.006$				$1.6 \pm 0.1$	$20 \pm 0$
	SSD	1e-02	$0.935 \pm 0.006$	$0.497 \pm 0.003$	$0.702 \pm 0.003$	$0.705 \pm 0.007$	$1.6 \pm 0.1$	$20\pm0$