

Rectifying Unlearning Efficacy and Privacy Evaluation: A New Inference Prospective

Nima Naderlou¹, Shenao Yan¹, Binghui Wang², Jie Fu³,
Wendy Hui Wang³, Weiran Liu⁴, Yuan Hong¹

¹University of Connecticut

²Illinois Institute of Technology

³Stevens Institute of Technology

⁴Alibaba Group

UConn

**ILLINOIS
TECH**



STEVENS
INSTITUTE of TECHNOLOGY
THE INNOVATION UNIVERSITY[®]

 **Alibaba**

USENIX Security 2025
Track 3: ML and AI privacy 2

It's 2025: Has Unlearning Already Won?

Every model has a "fast" unlearning fix.

A large and growing body of works have been introduced for inexact selective unlearning and, improvements are incremental

Empirical evaluations indicate that unlearning is approaching seamless “perfection”.



a) Unlearning request

It's 2025: Has Unlearning Already Won?

Failure of membership inference attack (MIA) → Better Forgetting [1]

Existing MIAs suggest that unlearning approximates **Retraining** (Gold standard)

Table 3: Performance of approximate unlearning methods (including both relabeling-free and relabeling-based methods) under random forget sets and worst-case forget sets on CIFAR-10 using ResNet-18 with forgetting ratio 10%. The result format follows Table 2. Additionally, a performance gap against **Retrain** is provided in (•). The metric *averaging (avg.) gap* is calculated by averaging the performance gaps measured in all metrics. Note that the better performance of an MU method corresponds to the smaller performance gap with Retrain.

Methods	Random Forget Set					Worst-Case Forget Set				
	UA	MIA	RA	TA	Avg. Gap	UA	MIA	RA	TA	Avg. Gap
Retrain	5.28 _{±0.33}	12.86 _{±0.61}	100.00 _{±0.00}	94.38 _{±0.15}	0.00	0.00 _{±0.00}	0.00 _{±0.00}	100.00 _{±0.00}	94.66 _{±0.09}	0.00
Relabeling-free										
FT	5.08 _{±0.39} (0.20)	10.96 _{±0.38} (1.90)	97.40 _{±0.52} (2.94)	91.02 _{±0.38} (3.36)	2.00	0.00 _{±0.00} (0.00)	0.02 _{±0.02} (0.02)	97.63 _{±0.46} (2.37)	91.58 _{±0.40} (3.08)	1.37
EU-k	2.34 _{±0.79} (2.94)	6.35 _{±0.89} (6.51)	97.52 _{±0.89} (2.48)	90.17 _{±0.88} (4.21)	4.04	0.68 _{±0.36} (0.68)	5.02 _{±4.42} (5.02)	97.17 _{±0.86} (2.83)	90.08 _{±0.70} (4.58)	3.28
CF-k	0.02 _{±0.02} (5.26)	0.76 _{±0.02} (12.10)	99.98 _{±0.00} (0.02)	94.45 _{±0.02} (0.07)	4.36	0.00 _{±0.00} (0.00)	0.00 _{±0.00} (0.00)	99.98 _{±0.01} (0.02)	94.34 _{±0.00} (0.32)	0.08
SCRUB	12.42 _{±0.92} (7.14)	22.63 _{±2.42} (9.57)	88.33 _{±0.76} (11.69)	83.15 _{±1.92} (11.20)	9.91	0.03 _{±0.03} (0.01)	0.04 _{±0.05} (0.04)	98.95 _{±0.33} (1.85)	92.75 _{±0.30} (1.88)	0.82
ℓ ₁ -sparse	4.34 _{±0.73} (0.94)							96.93 _{±0.73} (3.07)	90.96 _{±0.82} (3.70)	1.72

One-way MIA acc:
low MIA accuracy
gap < 3% with
“retraining” on
top unlearning
[2].

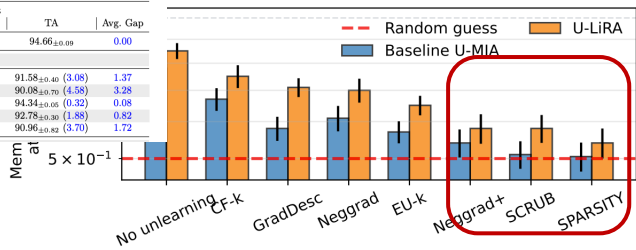
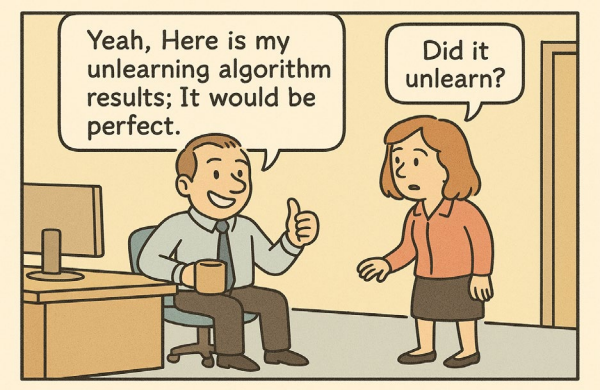


Figure 1 | Membership inference attack accuracy using a baseline attack and U-LiRA across different unlearning algorithms. Attack and unlearning algorithm descriptions are in Section 4. U-LiRA outperforms the baseline by a large margin across all unlearning algorithms because it creates per-example MIA decision rules.

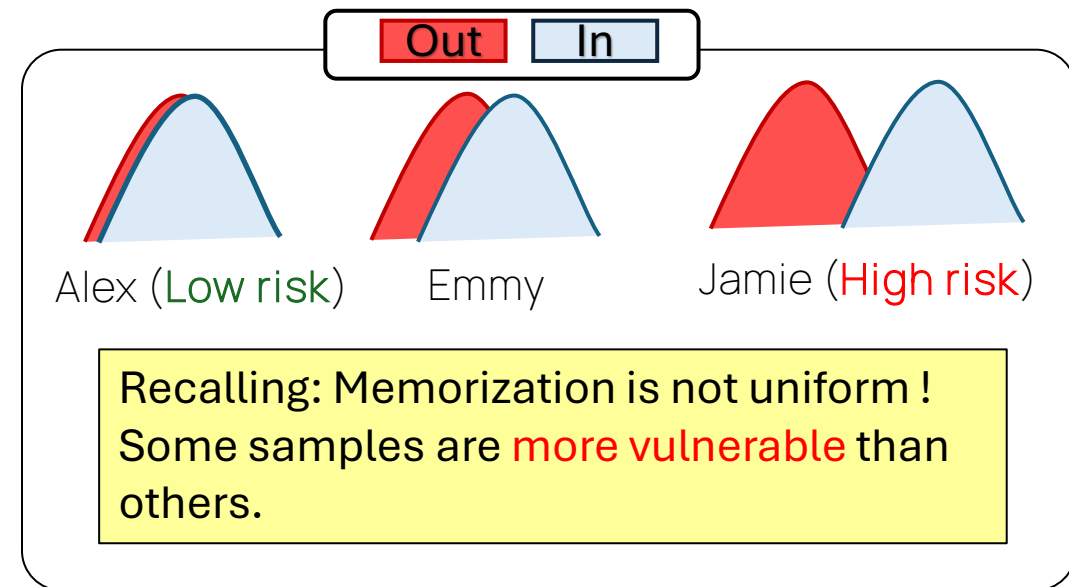
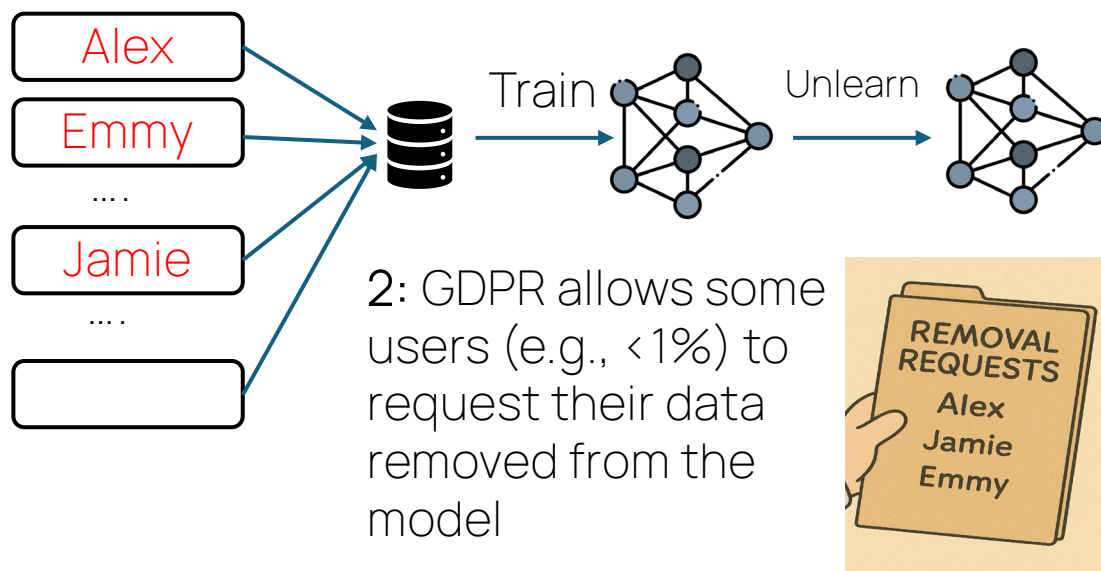
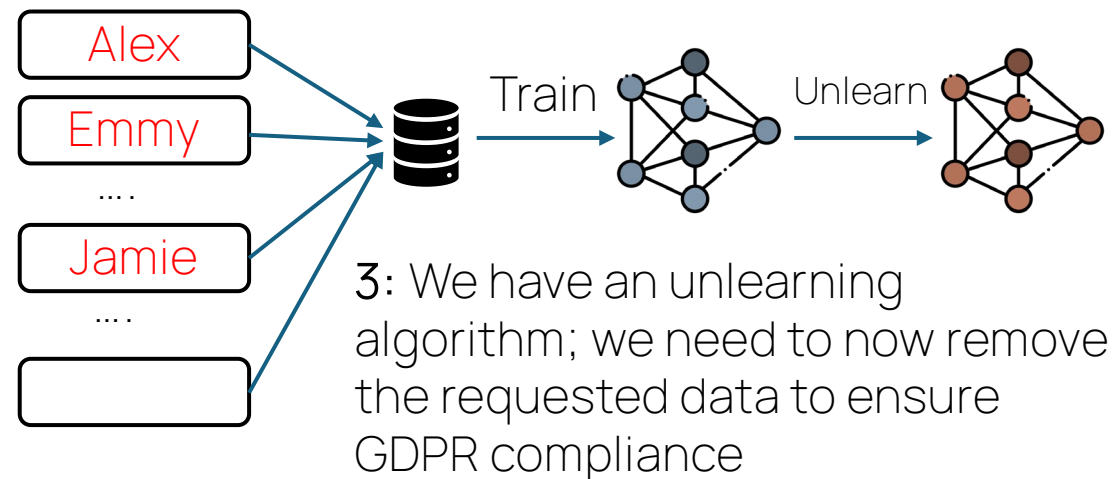
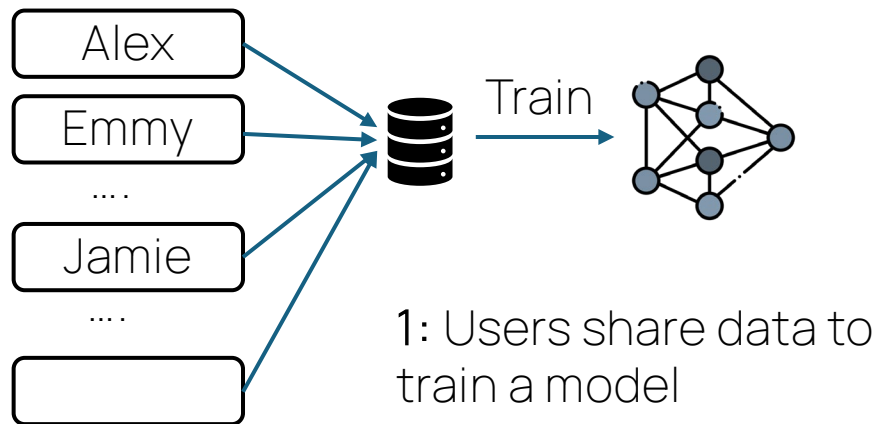


b) Using a fast inexact unlearning

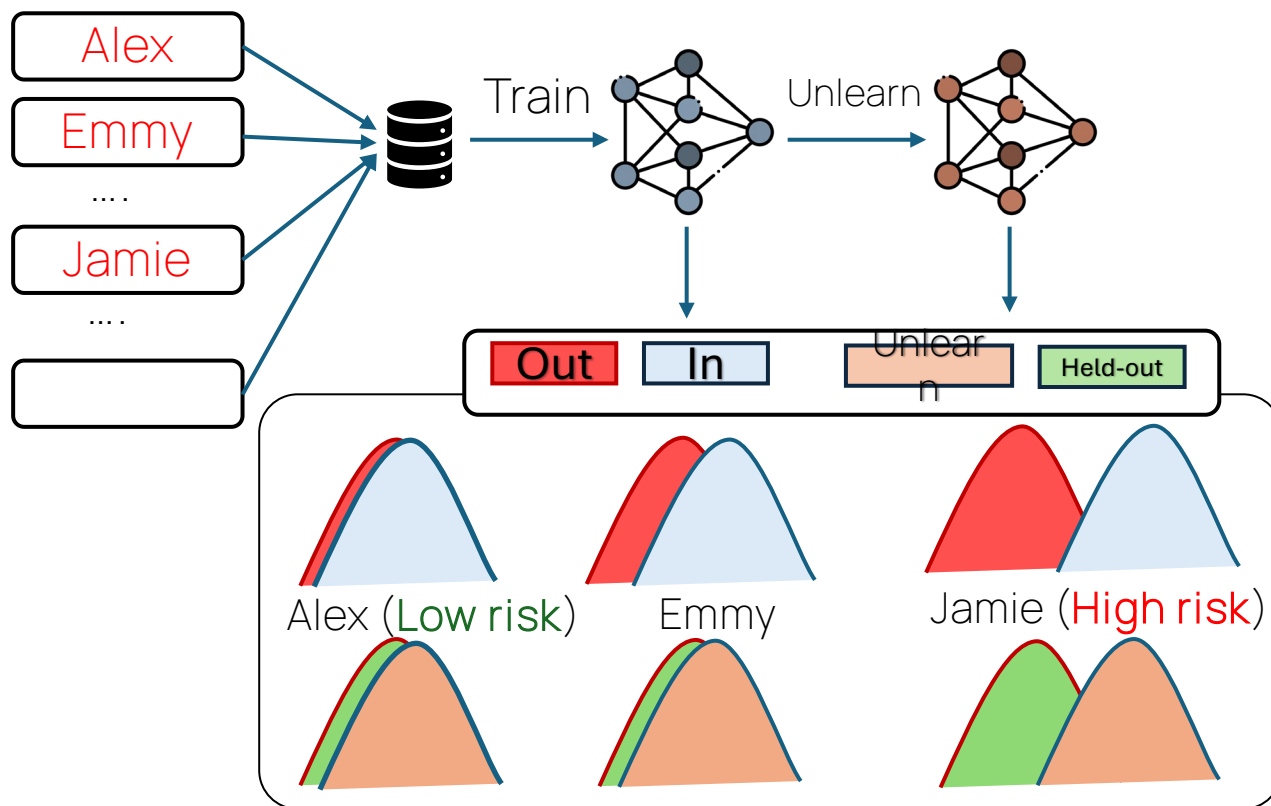
SOTA on privacy leakage:
MIA accuracy gap < 10%
on top unlearning [3].

[1] Jagielski, Matthew, et al. "Measuring forgetting of memorized training examples." In ICLR 2023.
[2] Fan, Chongyu, et al. "Challenging forgets: Unveiling the worst-case forget sets in machine unlearning." In ECCV 2024.
[3] Hayes, Jamie, et al. "Inexact unlearning needs more careful evaluations to avoid a false sense of privacy." In SaTML 2025.

Warmup up: Our motivation



What is missing today: Our motivation



If $Unlearn \approx Held-out$, privacy is protected.

"Privacy Leakage"

If $Unlearn \approx Out$, unlearning is effective.

"Efficacy" (Indistinguishability to Retraining)

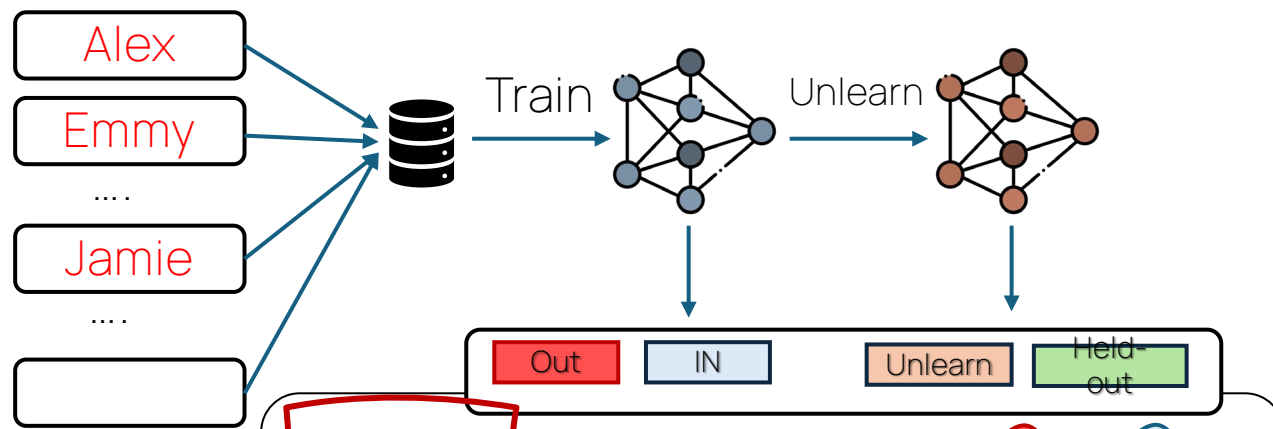
In: distribution of trained models where a sample is *member*

Out: distribution of trained models where sample is *non-member*

Unlearn: distribution of unlearned models where a sample is *unlearned*

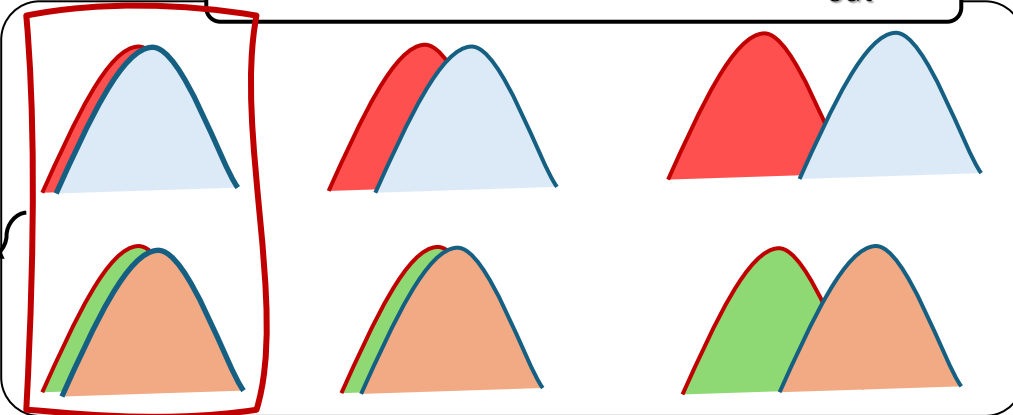
Held-out: distribution of unlearned models where sample is *non-member*

What is missing today: Our motivation



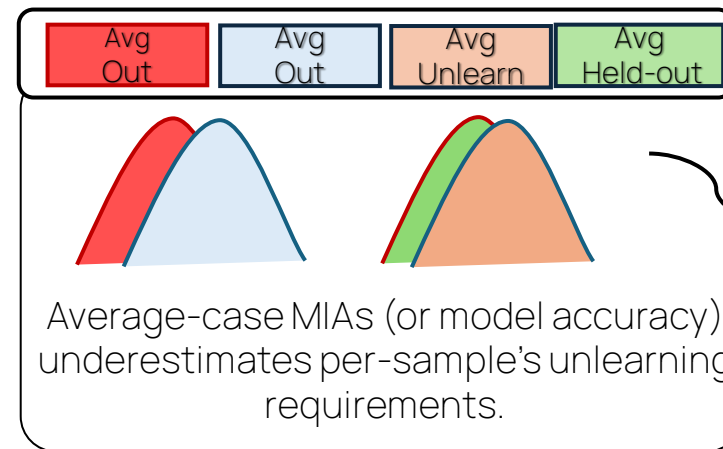
3. Many samples are like this; well-protected already.

“Let's not evaluate them”

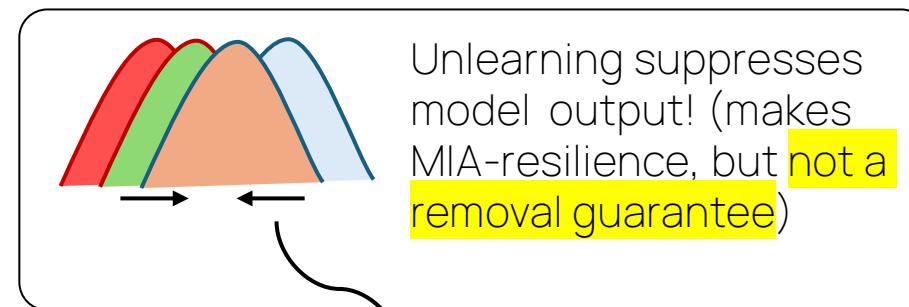


If $Unlearn \sim Held-out$ privacy is protected (Privacy leakage)

If $Unlearn \sim Out$ unlearning is effective (Efficacy) (Indistinguishability to Retraining)



1. “Better to be per-sample like [1]”



2. MIA resilience differs from unlearning guarantee! Need to find a way to measure efficacy

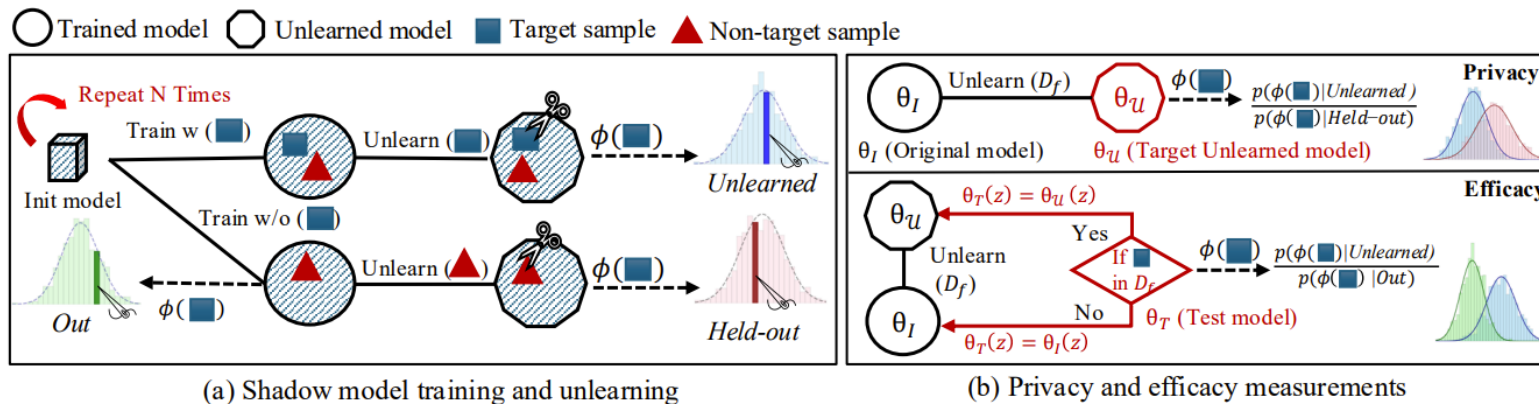
“MIA to detect if sample is unlearned or **retrained (Efficacy)**”

Our framework: RULI

1. We introduce an algorithm to train shadow models; get all distributions required per-sample

We tried optimizing parallelization our algorithm to minimize the shadow costs!

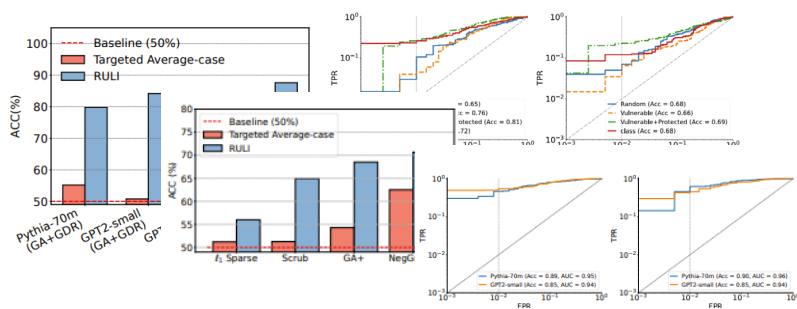
2. We introduce a hypothetical *Test model* to measure Efficacy ;
This calibrates output suppression impact.



3. We target vulnerable samples; went further and inject them as **canaries** to challenge unlearning.

Our Results

- ❖ We performed our attack on best-inexact unlearning baselines.
- ❖ We assume we can always find **best unlearning parameters** per unlearning request.
- ❖ Canary injection usually **leaks** more than purely unlearning vulnerable samples!
- ❖ We also tried similar experiments on: CIFAR-10, CIFAR-100 and. For generalizability, unlearning random 7-gram from GPT-2; similar trends exists!



Target data	Targeted average-case attack (Population attack)				RULI			
	AUC	ACC	TPR@ 1%FPR	TPR@ 5%FPR	AUC	ACC	TPR@ 1% FPR	TPR@ 5%FPR
ℓ_1 Sparse								
Vulnerable only	54.4%	55.1%	2.3%	5.2%	59.6%	56.0%	2.4%	12.4%
Vulnerable as canaries	55.3%	54.7%	0.8%	5.6%	62.6%	57.0%	6.3%	16.6%
Random	53.2%	52.8%	0.0%	2.4%	56%	54.4%	0.8%	6.4%
Scrub								
Vulnerable only	52.5%	52.4%	2.0%	5.4%	65.3%	61.5%	11.7%	23.9%
Vulnerable as canaries	56.0%	56.2%	1.0%	6.3%	69.5%	63.6%	10.9%	27.1%
Random	49.6%	49.8%	1.0%	2.8%	59.7%	57.0%	6.0%	14.0%

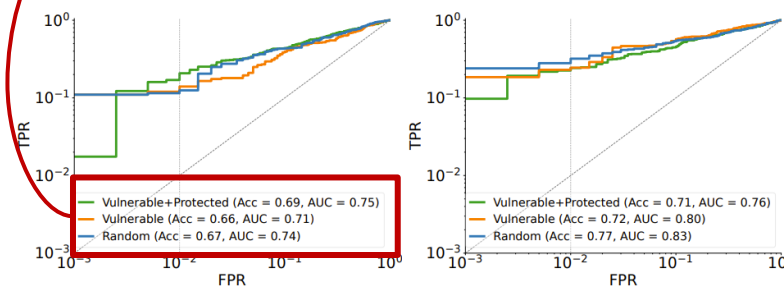
~12.6% higher
MIA success
x6.3 higher
privacy risk
than
retraining

~19.5% higher
MIA success;
x10.9 privacy
risk than
retraining

Tiny ImageNet unlearning; Swin-small model; unlearning <1% of the data.

500 samples: 250 Out and 250 Unlearned

Up to 69% MIA success
distinguishing
unlearned vs retrained



(a) ℓ_1 Sparse

(b) Scrub

An example of our results; Please find more results available in our paper!

- ❑ Our research *does not refute or critique* idea of inexact unlearning or existing evaluations. However, we claim:
- ❑ We need to *rectify & improve* our evaluation toward **stronger evaluations** to understand limitations of unlearning for privacy.

Thanks for your attention!

But we did not tell everything?!

Interested in more details about our design and experiments?

RULI is not perfect! Wondering about **RULI's limitations** and edge cases?

Let's discuss more in the following poster session

Or contact us via email: nima.naderlou@uconn.edu