

Motivation

Machine unlearning attracts many attentions. To evaluate it:

- *Membership Inference Attack* (MIA): most common choice
 \Rightarrow directly reflects **individual privacy risk**.

However, existing MIA-based evaluations are often not

1. **Well-calibrated** across different unlearning methods;
2. **Zero-grounded**: retraining is not always ranked highest;
3. **Comparable** across different attacks, yielding inconsistency.

Overview and Contributions

1. Formalize *unlearning sample inference game*, establishing a **novel unlearning evaluation metric** for data removal efficacy.
2. Demonstrate several **provable properties** of the proposed metric, dodging various pitfalls of existing MIA-based metrics.
3. Introduce a *SWAP* test for **efficient empirical** analysis.

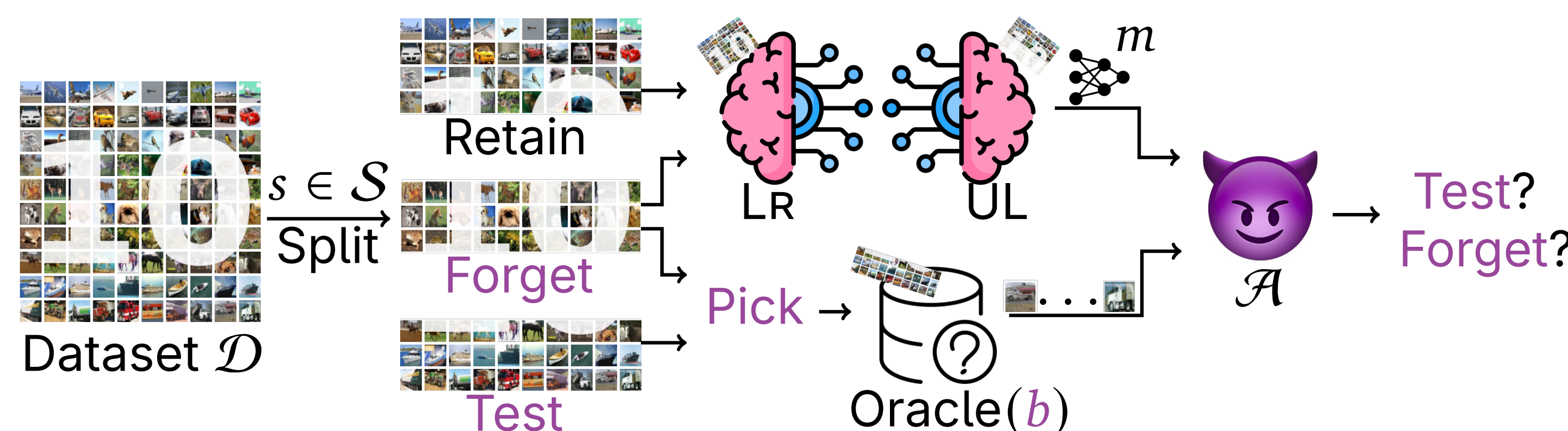
Machine Unlearning Evaluation as an Inference Game

We formulate unlearning as an **game** between

- UL: Unlearning algorithm (challenger), and
- \mathcal{A} : Membership-inference adversary \mathcal{A} .

Given a dataset \mathcal{D} , the *unlearning inference game* \mathcal{G} :

1. Split \mathcal{D} into retain, **forget**, and **test** sets, forming a split $s \in \mathcal{S}$.
2. *Random oracle* $O_s(b)$ with a **secret bit** $b \in \{0, 1\}$ is instantiated.
3. UL outputs unlearned model m , and \mathcal{A} attempts to infer b .



Question. How can we measure the performance of UL and \mathcal{A} ?

Cryptographic Advantage and Unlearning Quality

Intuitively, how well UL can **fool** \mathcal{A} measures the performance.

- This is a well-known concept in cryptography: **advantage**.

For our game \mathcal{G} , the advantage $\text{Adv}(\mathcal{A}, \text{UL})$ of \mathcal{A} against UL is

$$\frac{1}{|\mathcal{S}|} \left| \sum_{s \in \mathcal{S}} \Pr_{m \sim \mathbb{P}(\text{UL}, s)} (\mathcal{A}^O(m) = 1) - \sum_{s \in \mathcal{S}} \Pr_{m \sim \mathbb{P}(\text{UL}, s)} (\mathcal{A}^O(m) = 1) \right|.$$

Definition (Unlearning Quality). For any UL, its *Unlearning Quality* Q under a game \mathcal{G} is defined as

$$Q(\text{UL}) := 1 - \sup_{\mathcal{A}} \text{Adv}(\mathcal{A}, \text{UL}),$$

Theoretical Guarantees for Q

Theorem (Zero Grounding). For any adversary \mathcal{A} , we have $\text{Adv}(\mathcal{A}, \text{ReTRAIN}) = 0$. Hence, $Q(\text{ReTRAIN}) = 1$.

This guarantees that the retraining method is always the best.

Theorem (Calibrated Guarantees). Given an (ϵ, δ) -certified removal UL for some $\epsilon, \delta > 0$, for any \mathcal{A} against UL, we have

$$\text{Adv}(\mathcal{A}, \text{UL}) \leq 2 \cdot \left(1 - \frac{2 - 2\delta}{e^\epsilon + 1} \right) \Rightarrow Q(\text{UL}) \geq \frac{4 - 4\delta}{e^\epsilon + 1} - 1$$

Hence, Q **calibrates** with other known privacy metrics faithfully.

SWAP Test: Approximation Algorithm for Q

To efficiently evaluate the Q , we propose a *SWAP* test:

- Consider only **swapped splits** s, s' between forget and test set.
- Approximate $\text{Adv}(\mathcal{A}, \text{UL})$ by only few swap pairs.

Theorem (SWAP's Zero Grounding). For any \mathcal{A} and swap splits $s, s' \in \mathcal{S}$, $\text{Adv}_{\{s, s'\}}(\mathcal{A}, \text{ReTRAIN}) = 0$.

It turned out that SWAP is not only sufficient, but necessary.

Theorem (Blowup without SWAP). For two non-swapped splits $s_1, s_2 \in \mathcal{S}$, there exists \mathcal{A} such that $\text{Adv}_{\{s_1, s_2\}}(\mathcal{A}, \text{UL}) = 1$ for **any** UL. Particularly, $\text{Adv}_{\{s_1, s_2\}}(\mathcal{A}, \text{ReTRAIN}) = 1$.

Experimental Results: Model Trained with Different Privacy

Consider unlearning on models trained **with DP budgets** ϵ .

Unlearning Quality is negatively correlated with DP budget ϵ !

UL	ϵ				
	50	150	600	∞	
None	0.972 ⁺	0.960 ⁺	0.932 [*]	0.587 ⁺	
NegGrad	0.980 [*]	0.975 ⁺	0.953 ⁺	0.628 ⁺	
RetrFINAL	0.972 ⁺	0.964 ⁺	0.939 ⁺	0.576 ⁺	
FtFINAL	0.973 [*]	0.963 ⁺	0.939 ⁺	0.574 ⁺	
Fisher	0.973 [*]	0.967 ⁺	0.942 [*]	0.709 ⁺	
SalUN	0.979 [*]	0.972 ⁺	0.945 [*]	0.689 [*]	
SSD	0.996 [*]	0.988 [*]	0.981 ⁺	0.888 ⁺	
ReTRAIN	0.998 [*]	0.996 [*]	0.997 [*]	0.993 [*]	

(A) Q score versus DP budgets.

UL	ϵ				
	50	150	600	∞	
None	0.451 ⁺	0.433 ⁺	0.454 ⁺	0.380 ⁺	
NegGrad	0.476 ⁺	0.482 ⁺	0.466 ⁺	0.299 ⁺	
RetrFINAL	0.485 ⁺	0.485 ⁺	0.472 ⁺	0.248 ⁺	
FtFINAL	0.485 ⁺	0.485 ⁺	0.472 ⁺	0.247 ⁺	
Fisher	0.475 ⁺	0.484 ⁺	0.463 ⁺	0.325 ⁺	
SalUN	0.488 [*]	0.491 [*]	0.477 ⁺	0.268 [*]	
SSD	0.480 [*]	0.480 [*]	0.468 ⁺	0.244 [*]	
ReTRAIN	0.479 [*]	0.491 [*]	0.492 [*]	0.488 [*]	

(B) MIA score versus DP budgets.

⁺ indicates standard error of the mean is < 0.01 , and ^{*} for < 0.005 .

Next, we consider applying unlearning on different **dataset sizes**.

Unlearning Quality maintains a consistent ranking of UL!

UL	Dataset percentage (%)			
	0.1	0.4	0.8	1.0
RetrFINAL	0.340 \pm 0.017	0.586 \pm 0.015	0.621 \pm 0.014	0.634 \pm 0.025
FtFINAL	0.131 \pm 0.011	0.585 \pm 0.016	0.619 \pm 0.014	0.634 \pm 0.024
Fisher	0.751 \pm 0.024	0.679 \pm 0.005	0.734 \pm 0.006	0.791 \pm 0.020
NegGrad	0.124 \pm 0.010	0.564 \pm 0.018	0.603 \pm 0.014	0.656 \pm 0.035
SalUN	0.476 \pm 0.014	0.617 \pm 0.016	0.689 \pm 0.013	0.748 \pm 0.004
SSD	0.975 \pm 0.008	0.939 \pm 0.025	0.929 \pm 0.021	0.928 \pm 0.015
ReTRAIN	0.999 \pm 0.000	0.997 \pm 0.001	0.993 \pm 0.001	0.993 \pm 0.001

- **Well-calibrated**: Q not only calibrates under ϵ , but also other hyperparameters such as dataset percentage.
- **Zero-grounded**: For all settings, $Q(\text{ReTRAIN}) \approx 1$.
- **Comparable**: While MIA score is inconsistent, Q unifies it.

Next Step

1. Efficient adaptation to foundation models unlearning?
2. More complicated unlearning scenarios, such as non-i.i.d. unlearning and feature unlearning?