



**ICLR**

# Membership Inference Attacks Against Fine-tuned Diffusion Language Models

Yuetian Chen<sup>1</sup>, Kaiyuan Zhang<sup>1</sup>, Yuntao Du<sup>1</sup>, Edoardo Stoppa<sup>1</sup>, Charles Fleming<sup>2</sup>, Ashish Kundu<sup>2</sup>,  
Bruno Ribeiro<sup>1</sup>, Ninghui Li<sup>1</sup>

<sup>1</sup> Department of Computer Science, Purdue University;

<sup>2</sup> Cisco Research

Presenter: Yuetian Chen

# Context: Privacy Risks in Diffusion Language Models

*Do the unique paradigms of DLMs introduce novel vulnerabilities exploitable by MIAs?*

- **A Distinct Paradigm: Diffusion Language Models (DLMs)**

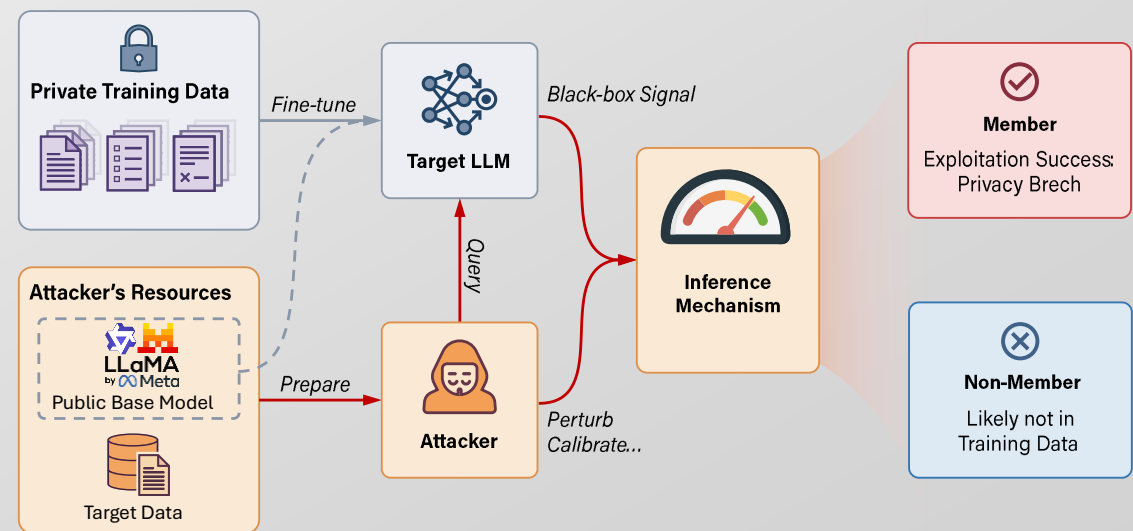
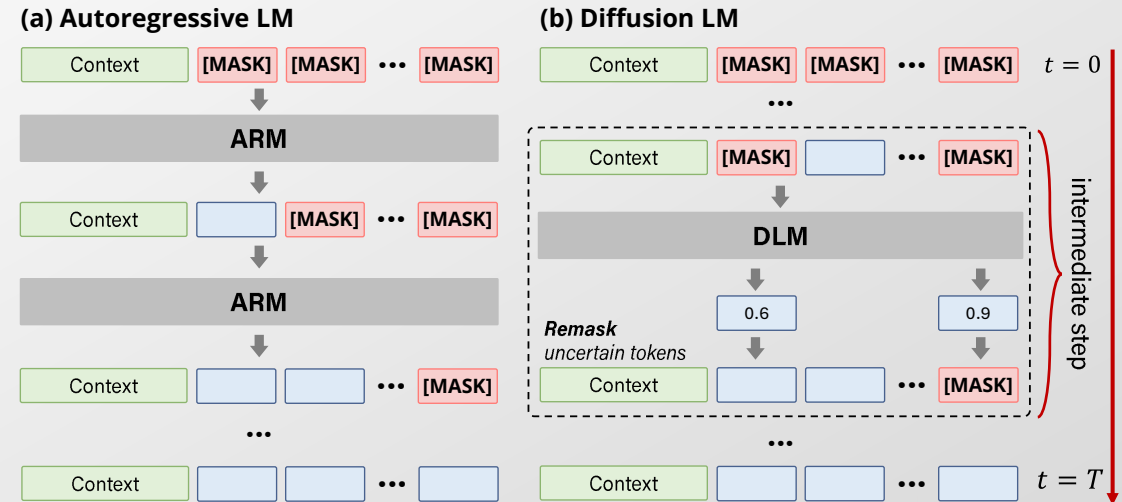
- DLMs are emerging as a compelling alternative to Autoregressive Models (ARMs)
- **Key Difference:** Instead of fixed sequential prediction, DLMs use bidirectional context to iteratively reconstruct masked tokens

- **The Threat: Membership Inference Attacks (MIA)**

- **Goal:** Detect if a specific text sample (e.g., private user data) was used to fine-tune the target model
- **Threat Model:** We assume a standard Grey-box adversary with query access (logits) to the fine-tuned model and access to a reference model

- **The Opportunity: Exploiting "Extra" Signals**

- The masking mechanism allows attackers to probe multiple, arbitrary configurations of the same text in DLMs
- **Hypothesis:** This flexibility exposes richer, multi-view membership signals that existing attacks fail to capture

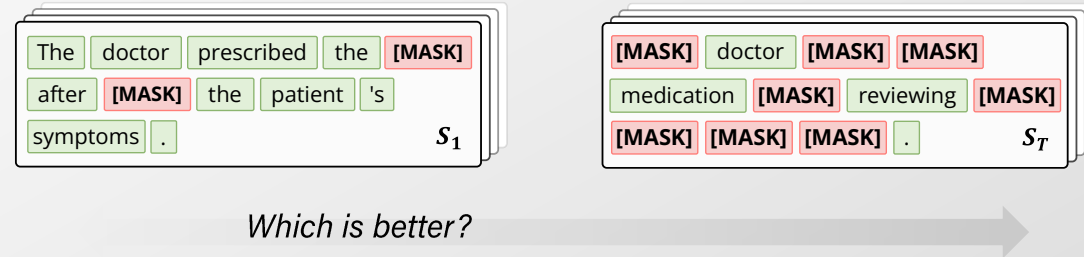


# The Challenges of Attacking DLMs

*How we mask? & How we score?*

## • The Masking Dilemma

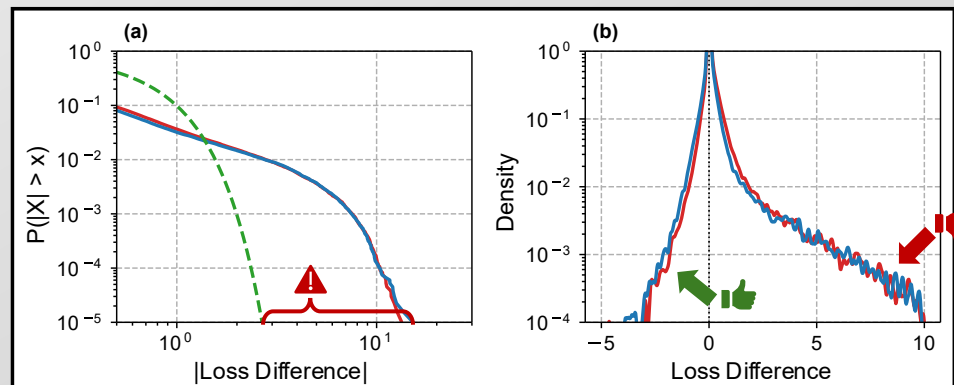
- A random mask configuration will likely miss the specific token pattern the model memorized
- You mostly get weak signals dominated by pure noise



## • The Problem with Loss Averaging

- Long tail noise make loss averaging unreliable
- **Domain tokens are "loud"**: Extreme token loss spikes usually come from ubiquitous "domain adaptation" terms (e.g., "Biography")
- **Contamination**: These loss outliers contain **little membership info**, yet their sheer magnitude completely dominates the average score

Most Frequent Token	Loss Reduction after Fine-tune	Membership Signal Strength	Count in Corpus
' Biography '	1.040	0.295	50
' Background '	0.671	0.027	25
' Life '	0.631	0.299	21
' List '	0.544	0.499	39
' Description '	0.399	0.474	23
' Joseph '	0.393	0.022	23



# Methodology: SAMA

## Subset-Aggregated Membership Attack

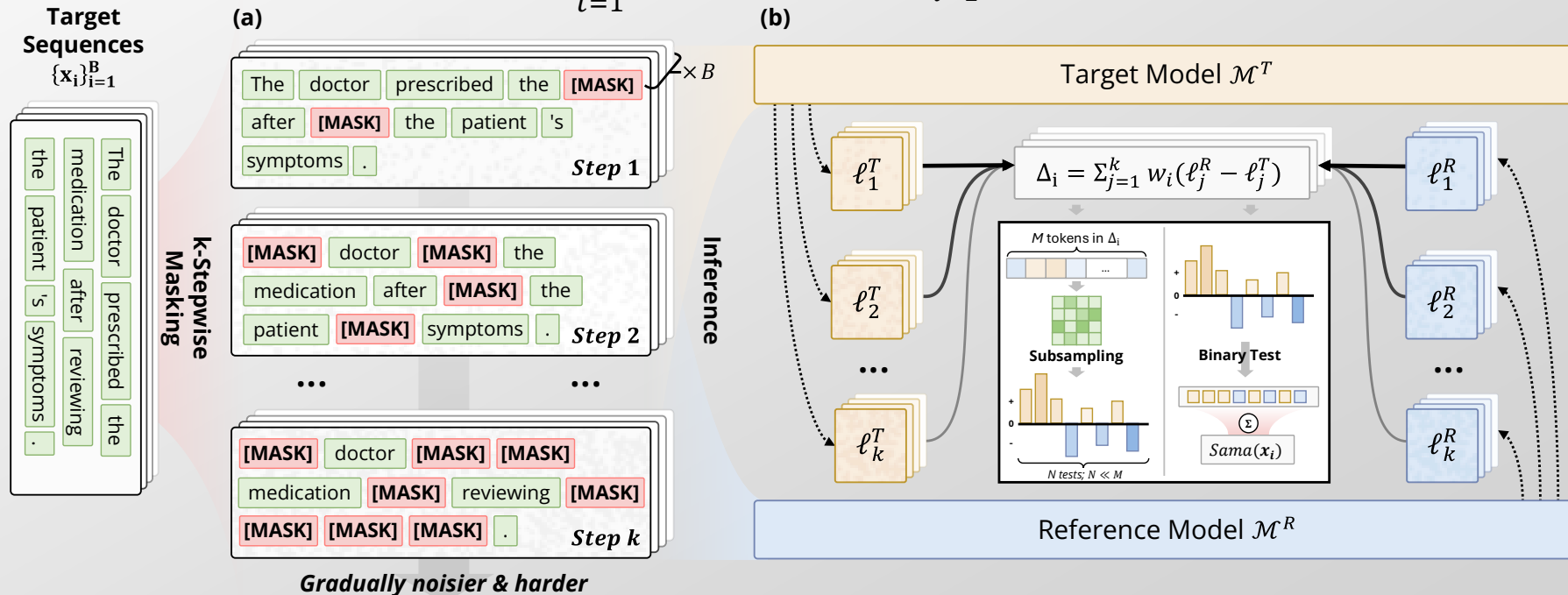
- (a) Solving The Masking Dilemma - Progressive Evidence Collection

- Memorization is a moving target.** We triangulate memorization across multiple densities.
- Multi-Scale Probing:** Systematically sweep masking density  $\alpha_t$  over  $T$  steps to capture patterns at every context level:

$$\alpha_t = \alpha_{min} + \frac{t - 1}{T - 1} (\alpha_{max} - \alpha_{min})$$

- Adaptive Weighting:** Inverse-step weighting automatically prioritizes the cleaner, high-fidelity signals found in sparse masks:

$$SAMA(x) = \sum_{t=1}^T w_t \hat{\beta}_t(x), \quad w_t = \frac{1/t}{\sum_{i=1}^T 1/i}$$



# Methodology: SAMA

## *Subset-Aggregated Membership Attack*

- **(b) Solving Long-tail Noise - Robust Subset Aggregation**

- Sign-based estimators neutralize infinite-variance domain noise, isolating sparse memorization signals.
- **The Statistical Dilemma:** Standard Monte Carlo estimates of the expected loss gap,  $\mathbb{E}_{\mathcal{S}}[\Delta_{DF}(x; \mathcal{S})]$ , are easily hijacked by heavy-tailed domain adaptation artifacts (excess kurtosis  $> 80$ ).
- **Localized Subsampling:** We isolate extreme outliers by drawing  $N$  random subsets  $\mathcal{U}^n \subset \mathcal{S}_t$  to compute localized loss gaps:

$$\Delta_{DF}^n(x; \mathcal{S}_t) = \frac{1}{m} \sum_{i \in \mathcal{U}^n} [l_i^R(x, \mathcal{S}_t) - l_i^T(x, \mathcal{S}_t)]$$

- **Sign-Based Transformation:** Leveraging Hodges-Lehmann robust statistics, we map the continuous gap to a binary indicator, discarding magnitude entirely

$$B^n(x) = \mathbf{1}[\Delta_{DF}^n(x; \mathcal{S}_t) > 0]$$

- **The Theoretical Guarantee:** For non-members, symmetric noise ensures  $\mathbb{E}[B^n] = 0.5$  regardless of infinite variance. For members, consistent memorization breaks this symmetry, dominating the aggregated vote  $\hat{\beta}_t(x) = \frac{1}{N} \sum B^n(x)$ .

# Experimental Results

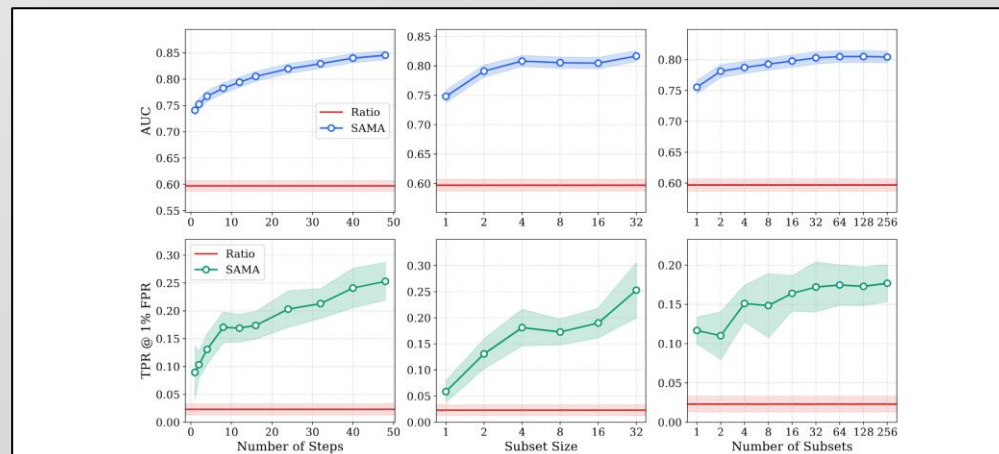
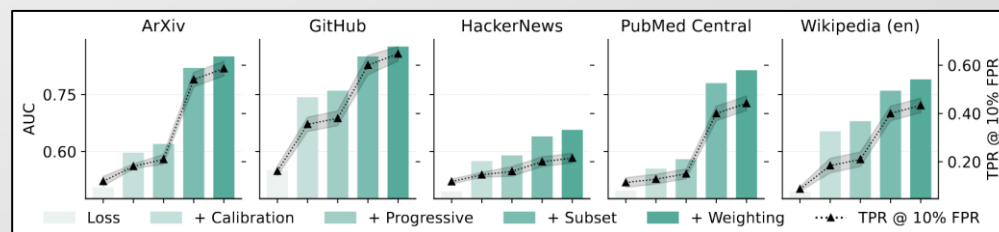
Table 1: MIA performance (AUC, TPR@10%FPR, TPR@1%FPR, TPR@0.1%FPR) across datasets. Each cell shows the mean with std. dev. as a subscript. The best results are highlighted.

MIAs	ArXiv				GitHub				HackerNews			
	AUC	T@10%	T@1%	T@0.1%	AUC	T@10%	T@1%	T@0.1%	AUC	T@10%	T@1%	T@0.1%
Loss	0.506 $\pm$ .01	0.119 $\pm$ .02	0.010 $\pm$ .01	0.000 $\pm$ .00	0.551 $\pm$ .01	0.161 $\pm$ .02	0.036 $\pm$ .01	0.007 $\pm$ .01	0.495 $\pm$ .01	0.118 $\pm$ .01	0.010 $\pm$ .01	0.000 $\pm$ .00
ZLIB	0.490 $\pm$ .01	0.119 $\pm$ .02	0.012 $\pm$ .00	0.000 $\pm$ .00	0.561 $\pm$ .01	0.205 $\pm$ .02	0.045 $\pm$ .01	0.007 $\pm$ .00	0.486 $\pm$ .01	0.083 $\pm$ .01	0.009 $\pm$ .01	0.001 $\pm$ .00
Lowercase	0.515 $\pm$ .01	0.103 $\pm$ .01	0.013 $\pm$ .00	0.001 $\pm$ .00	0.579 $\pm$ .01	0.178 $\pm$ .02	0.044 $\pm$ .01	0.008 $\pm$ .00	0.483 $\pm$ .01	0.074 $\pm$ .01	0.007 $\pm$ .00	0.001 $\pm$ .00
Min-K%	0.488 $\pm$ .01	0.102 $\pm$ .01	0.012 $\pm$ .00	0.001 $\pm$ .00	0.530 $\pm$ .01	0.171 $\pm$ .02	0.039 $\pm$ .01	0.007 $\pm$ .01	0.492 $\pm$ .01	0.108 $\pm$ .01	0.013 $\pm$ .01	0.001 $\pm$ .00
Min-K%++	0.485 $\pm$ .01	0.095 $\pm$ .01	0.006 $\pm$ .00	0.000 $\pm$ .00	0.496 $\pm$ .01	0.117 $\pm$ .01	0.016 $\pm$ .01	0.003 $\pm$ .00	0.486 $\pm$ .01	0.100 $\pm$ .02	0.008 $\pm$ .00	0.002 $\pm$ .00
BoWs	0.519 $\pm$ .01	0.107 $\pm$ .01	0.011 $\pm$ .00	0.000 $\pm$ .00	0.656 $\pm$ .02	0.306 $\pm$ .02	0.154 $\pm$ .02	0.059 $\pm$ .05	0.527 $\pm$ .01	0.128 $\pm$ .01	0.009 $\pm$ .00	0.001 $\pm$ .00
ReCall	0.501 $\pm$ .01	0.132 $\pm$ .02	0.007 $\pm$ .00	0.000 $\pm$ .00	0.562 $\pm$ .01	0.187 $\pm$ .01	0.037 $\pm$ .01	0.007 $\pm$ .00	0.494 $\pm$ .01	0.090 $\pm$ .01	0.010 $\pm$ .01	0.000 $\pm$ .00
CON-Recall	0.500 $\pm$ .02	0.101 $\pm$ .02	0.011 $\pm$ .00	0.000 $\pm$ .00	0.549 $\pm$ .01	0.168 $\pm$ .01	0.027 $\pm$ .01	0.005 $\pm$ .00	0.501 $\pm$ .02	0.098 $\pm$ .01	0.015 $\pm$ .01	0.000 $\pm$ .00
Neighbor	0.506 $\pm$ .01	0.098 $\pm$ .01	0.012 $\pm$ .01	0.001 $\pm$ .00	0.478 $\pm$ .01	0.072 $\pm$ .01	0.008 $\pm$ .01	0.000 $\pm$ .00	0.520 $\pm$ .01	0.123 $\pm$ .01	0.009 $\pm$ .00	0.001 $\pm$ .00
Ratio	0.597 $\pm$ .01	0.181 $\pm$ .01	0.023 $\pm$ .01	0.001 $\pm$ .00	0.743 $\pm$ .01	0.355 $\pm$ .03	0.081 $\pm$ .02	0.017 $\pm$ .01	0.575 $\pm$ .02	0.146 $\pm$ .01	0.013 $\pm$ .01	0.000 $\pm$ .00
SecMI	0.520 $\pm$ .01	0.096 $\pm$ .02	0.006 $\pm$ .00	0.001 $\pm$ .00	0.604 $\pm$ .01	0.190 $\pm$ .01	0.044 $\pm$ .01	0.019 $\pm$ .01	0.523 $\pm$ .02	0.125 $\pm$ .02	0.013 $\pm$ .00	0.001 $\pm$ .00
PIA	0.525 $\pm$ .01	0.099 $\pm$ .01	0.011 $\pm$ .01	0.000 $\pm$ .00	0.571 $\pm$ .01	0.131 $\pm$ .01	0.012 $\pm$ .00	0.001 $\pm$ .00	0.494 $\pm$ .01	0.086 $\pm$ .01	0.014 $\pm$ .00	0.000 $\pm$ .00
<b>SAMA (Ours)</b>	<b>0.850<math>\pm</math>.01</b>	<b>0.586<math>\pm</math>.03</b>	<b>0.178<math>\pm</math>.03</b>	<b>0.014<math>\pm</math>.01</b>	<b>0.876<math>\pm</math>.01</b>	<b>0.647<math>\pm</math>.03</b>	<b>0.259<math>\pm</math>.05</b>	<b>0.075<math>\pm</math>.05</b>	<b>0.657<math>\pm</math>.01</b>	<b>0.215<math>\pm</math>.02</b>	<b>0.027<math>\pm</math>.01</b>	<b>0.003<math>\pm</math>.00</b>

MIAs	PubMed Central				Wikipedia (en)				Pile CC			
	AUC	T@10%	T@1%	T@0.1%	AUC	T@10%	T@1%	T@0.1%	AUC	T@10%	T@1%	T@0.1%
Loss	0.498 $\pm$ .01	0.114 $\pm$ .02	0.016 $\pm$ .01	0.001 $\pm$ .00	0.495 $\pm$ .01	0.087 $\pm$ .00	0.010 $\pm$ .00	0.003 $\pm$ .00	0.502 $\pm$ .01	0.105 $\pm$ .01	0.012 $\pm$ .00	0.002 $\pm$ .00
ZLIB	0.488 $\pm$ .01	0.096 $\pm$ .02	0.005 $\pm$ .00	0.001 $\pm$ .00	0.495 $\pm$ .01	0.093 $\pm$ .01	0.008 $\pm$ .00	0.000 $\pm$ .00	0.491 $\pm$ .01	0.095 $\pm$ .01	0.007 $\pm$ .00	0.001 $\pm$ .00
Lowercase	0.502 $\pm$ .01	0.096 $\pm$ .01	0.002 $\pm$ .00	0.000 $\pm$ .00	0.535 $\pm$ .01	0.107 $\pm$ .02	0.013 $\pm$ .00	0.001 $\pm$ .00	0.518 $\pm$ .01	0.102 $\pm$ .01	0.009 $\pm$ .00	0.001 $\pm$ .00
Min-K%	0.500 $\pm$ .01	0.119 $\pm$ .01	0.008 $\pm$ .00	0.002 $\pm$ .00	0.482 $\pm$ .01	0.070 $\pm$ .01	0.008 $\pm$ .00	0.000 $\pm$ .00	0.491 $\pm$ .01	0.095 $\pm$ .01	0.008 $\pm$ .00	0.001 $\pm$ .00
Min-K%++	0.494 $\pm$ .01	0.109 $\pm$ .01	0.011 $\pm$ .00	0.000 $\pm$ .00	0.488 $\pm$ .01	0.118 $\pm$ .02	0.004 $\pm$ .00	0.000 $\pm$ .00	0.491 $\pm$ .01	0.113 $\pm$ .01	0.007 $\pm$ .00	0.001 $\pm$ .00
BoWs	0.489 $\pm$ .01	0.100 $\pm$ .01	0.002 $\pm$ .00	0.000 $\pm$ .00	0.471 $\pm$ .01	0.084 $\pm$ .02	0.003 $\pm$ .00	0.001 $\pm$ .00	0.480 $\pm$ .01	0.092 $\pm$ .01	0.003 $\pm$ .00	0.001 $\pm$ .00
ReCall	0.495 $\pm$ .01	0.088 $\pm$ .01	0.007 $\pm$ .00	0.000 $\pm$ .00	0.506 $\pm$ .01	0.103 $\pm$ .01	0.010 $\pm$ .01	0.000 $\pm$ .00	0.500 $\pm$ .01	0.096 $\pm$ .01	0.009 $\pm$ .00	0.000 $\pm$ .00
CON-Recall	0.498 $\pm$ .02	0.119 $\pm$ .02	0.009 $\pm$ .00	0.000 $\pm$ .00	0.498 $\pm$ .02	0.092 $\pm$ .01	0.008 $\pm$ .00	0.000 $\pm$ .00	0.498 $\pm$ .01	0.106 $\pm$ .01	0.009 $\pm$ .00	0.000 $\pm$ .00
Neighbor	0.506 $\pm$ .01	0.091 $\pm$ .01	0.011 $\pm$ .01	0.000 $\pm$ .00	0.504 $\pm$ .01	0.101 $\pm$ .01	0.009 $\pm$ .00	0.001 $\pm$ .00	0.505 $\pm$ .01	0.096 $\pm$ .01	0.010 $\pm$ .00	0.001 $\pm$ .00
Ratio	0.555 $\pm$ .01	0.128 $\pm$ .02	0.018 $\pm$ .01	0.003 $\pm$ .01	0.653 $\pm$ .01	0.184 $\pm$ .03	0.011 $\pm$ .01	0.000 $\pm$ .00	0.604 $\pm$ .01	0.156 $\pm$ .02	0.015 $\pm$ .01	0.002 $\pm$ .00
SecMI	0.510 $\pm$ .01	0.109 $\pm$ .01	0.009 $\pm$ .00	0.008 $\pm$ .00	0.522 $\pm$ .01	0.101 $\pm$ .01	0.015 $\pm$ .00	0.009 $\pm$ .00	0.515 $\pm$ .01	0.108 $\pm$ .01	0.010 $\pm$ .00	0.001 $\pm$ .00
PIA	0.496 $\pm$ .01	0.102 $\pm$ .01	0.005 $\pm$ .00	0.000 $\pm$ .00	0.522 $\pm$ .01	0.120 $\pm$ .01	0.011 $\pm$ .00	0.001 $\pm$ .00	0.509 $\pm$ .01	0.103 $\pm$ .01	0.009 $\pm$ .00	0.000 $\pm$ .00
<b>SAMA (Ours)</b>	<b>0.814<math>\pm</math>.01</b>	<b>0.442<math>\pm</math>.03</b>	<b>0.132<math>\pm</math>.03</b>	<b>0.011<math>\pm</math>.01</b>	<b>0.790<math>\pm</math>.01</b>	<b>0.433<math>\pm</math>.03</b>	<b>0.136<math>\pm</math>.01</b>	<b>0.008<math>\pm</math>.02</b>	<b>0.778<math>\pm</math>.01</b>	<b>0.408<math>\pm</math>.03</b>	<b>0.115<math>\pm</math>.02</b>	<b>0.009<math>\pm</math>.01</b>

## TL;DR

- **SOTA:** Outperforms 12 ARM and diffusion baselines.
- **Ablation:** Every SAMA component is necessary.
- **Robustness:** SAMA is highly stable across hyperparameter changes.





**ICLR**

# Membership Inference Attacks Against Fine-tuned Diffusion Language Models

Yuetian Chen<sup>1</sup>, Kaiyuan Zhang<sup>1</sup>, Yuntao Du<sup>1</sup>, Edoardo Stoppa<sup>1</sup>, Charles Fleming<sup>2</sup>, Ashish Kundu<sup>2</sup>,  
Bruno Ribeiro<sup>1</sup>, Ninghui Li<sup>1</sup>

<sup>1</sup> Department of Computer Science, Purdue University;

<sup>2</sup> Cisco Research

Thank you!



*Paper*



*Code*



*Blog*