

UniAud: A Unified Auditing Framework for High Auditing Power and Utility with One Training Run

Ruixuan Liu
Emory University
ruixuan.liu2@emory.edu

Li Xiong
Emory University
lxiong@emory.edu

Abstract—Differentially private (DP) optimization has been widely adopted as a standard approach to provide rigorous privacy guarantees for training datasets. DP auditing verifies whether a model trained with DP optimization satisfies its claimed privacy level by estimating empirical privacy lower bounds through hypothesis testing. Recent $O(1)$ frameworks improve auditing efficiency by checking the membership status of multiple audit samples in a single run, rather than checking individual samples across multiple runs. However, we reveal that there is no free lunch for this improved efficiency: data dependency and an implicit conflict between auditing and utility impair the tightness of the auditing results. Addressing these challenges, our key insights include reducing data dependency through uncorrelated data and resolving the auditing-utility conflict by decoupling the criteria for effective auditing and separating objectives for utility and auditing. We first propose a unified framework, UniAud, for data-independent auditing that maximizes auditing power through a novel uncorrelated canary construction and a self-comparison framework. We then extend this framework as UniAud++ for data-dependent auditing, optimizing the auditing and utility trade-off through multi-task learning with separate objectives for auditing and training. Experimental results validate that our black-box $O(1)$ framework matches the state-of-the-art auditing results of $O(T)$ auditing with thousands of runs, demonstrating the best efficiency-auditing trade-off across vision and language tasks. Additionally, our framework provides meaningful auditing with only slight utility degradation compared to standard DP training, showing the optimal utility-auditing trade-off and the benefit of requiring no extra training for auditing.

I. INTRODUCTION

As machine learning models gain widespread adoption in sensitive domains, data privacy has emerged as a critical concern. A compelling body of evidence demonstrates that trained models can inadvertently leak private information about their training data—through model gradients [1], loss values [2], or output predictions [3]. These findings underscore the urgent need for rigorous privacy guarantees throughout the learning process.

Differentially private (DP) training, such as DP-SGD [4], has become a standard approach for privacy-preserving machine learning. According to the DP definition [5], the probability that an adversary can distinguish between two models trained on datasets differing by a single data point is bounded by the privacy parameters (ϵ, δ) . This inherent guarantee makes DP a natural defense against membership inference attacks (MIAs) [2], [6], which aim to determine whether a particular individual’s data was included in the training set.

Moreover, DP has been adopted as a foundational component in mitigating other privacy attacks [3].

While DP training has been widely adopted, the strength of its guarantees critically depends on correct implementation. To achieve the claimed privacy level, DP training requires random batch sampling, per-sample gradient clipping, and gradient perturbation, where any flawed implementation can result in higher empirical privacy risk [7]. In contrast to the upper bound provided by DP, DP auditing estimates a lower bound on the actual privacy risk. A tight audit implies that the empirical lower bound closely approaches the theoretical upper bound. From the analyzer’s perspective, DP auditing serves to assess the tightness of the analytical bounds. From the practitioner’s perspective, it helps uncover implementation flaws, such as incorrect noise injection or faulty gradient clipping, in DP training systems.

Essentially, DP auditing works by formulating the process as a hypothesis test that determines whether to reject the null hypothesis: that the target model satisfies its claimed privacy level at a given confidence level. It relies on empirical observations by guessing the membership of an audit sample, in alignment with DP guarantees. The original auditing algorithms [8], [9], [10] have $O(T)$ complexity, as they require repeatedly training models on neighboring datasets that differ by one sample, for T observations. A recent $O(1)$ method [11] improves efficiency by randomly including or excluding multiple audit samples for obtaining multiple observations with a single training run.

Challenges in Black-Box $O(1)$ Auditing. Despite recent progress, a critical trade-off persists: $O(T)$ -methods [8], [9] achieve tighter empirical lower bounds (i.e., smaller gaps to the theoretical ϵ) by leveraging a greater number of observations, corresponding to T training runs. In contrast, $O(1)$ -methods [11], [12], [13] sacrifice auditing power for scalability, especially when the number of audit samples m in a single training run increases. This gap is further exacerbated in black-box settings, where auditors have access only to model outputs (e.g., APIs or prediction interfaces)—a realistic scenario in third-party audits or regulatory compliance. For example, black-box $O(1)$ auditing is significantly looser than its white-box counterpart [11], [12], potentially understating actual privacy risks. Nevertheless, black-box access is often necessary [9], [14] due to practical constraints in commercial deployments or limited system visibility available to auditors.

It remains an open problem to bridge the gap between efficiency and auditing power in black-box $O(1)$ auditing. This limitation undermines trust in DP deployments where transparency is limited—precisely the scenarios in which reliable auditing for the training algorithm is most critical.

We note that the challenge of the auditing-efficiency trade-off in $O(1)$ frameworks stems from two key factors. On the one hand, auditing power is weakened by data dependencies within real-world audit datasets. Correlated audit samples influence each other’s membership inference decisions—not only through interactions during training, but also because existing $O(1)$ auditing algorithms rely on relative ranking to predict membership. On the other hand, current $O(1)$ frameworks implicitly introduce an auditing-utility trade-off, complicating the auditing-efficiency dilemma into a trilemma. Existing methods [15] attempt to achieve both meaningful auditing results and strong model utility, but reconciling these two objectives is inherently difficult.

Key Insights. Our goal is to improve the auditing power of the efficient $O(1)$ framework by addressing the aforementioned challenges. First, we observe that data dependency manifests not only during training but also in the membership inference phase of existing $O(1)$ auditing algorithms. Thus, our key intuition is to minimize data dependency across both the training and inference stages. Second, we find that current DP auditing goals lack clear definitions when simultaneously considering the three aspects of auditing power, efficiency, and utility. Our key insight is to disentangle these factors—particularly by decoupling utility and auditing—while preserving the computational efficiency of $O(1)$ methods.

Our Solutions. To minimize data dependency during training, we propose a unified random pair-matching task that encodes the membership of each sample independently. This design promotes greater independence in the influence of each audit sample on the model. To reduce data dependency during the membership inference phase, we introduce a novel self-comparison-based inference method that removes reliance on relative ranking across samples, thereby making the MIA decision for each audit sample independent. To disentangle auditing and utility goals under efficient $O(1)$ computation, we explicitly distinguish between auditing objectives in data-independent and data-dependent settings—a distinction that has not been addressed in prior work. For data-independent auditing, which aims to evaluate the DP implementation independent of the training data or model, we design a separate auditing dataset and maximize auditing capability under $O(1)$ efficiency. For data-dependent auditing, which targets estimating the empirical privacy risk lower bound for a specific training dataset and model, we consider the triplet trade-offs and resolve the auditing-utility conflict through a decoupled training objective.

We summarize our contribution as follows:

- We deeply analyze the challenge in existing $O(1)$ framework, and reveal two fundamental insights for a unified black-box $O(1)$ auditing solution: reducing data dependency and decoupling utility from auditing. Thus, we

clarify data-independent and data dependent auditing goals and set up for the unified auditing framework, which is also general for different tasks and input formats (e.g., images and text sequences).

- To improve the auditing-efficiency trade-off in the data-independent setting, we propose the basic framework UniAud, which addresses the key challenge of data dependency through a novel pair-matching auditing task with synthetic canaries and self-comparison inference.
- To optimize the trade-off among auditing, efficiency, and utility in data-dependent auditing, we extend UniAud to UniAud++ through a novel design that decouples the training objectives of membership encoding (for auditing) and the main task (for utility).
- Through extensive evaluation on both image classification and language modeling tasks, we show that UniAud significantly improves the auditing performance of $O(1)$ methods in the data-independent setting, matching the state-of-the-art $O(T)$ auditing performance, while maintaining $O(1)$ efficiency and saving $\times 10^3$ of training runs compared to $O(T)$ methods; UniAud++ achieves the best overall trade-off compared to existing methods in the data-dependent setting, demonstrating the potential of integrating auditing along with normal training without extra runs.

II. RELATED WORKS

A. Theoretical Privacy Auditing

Privacy auditing [11], [8], [9], [10] has emerged as a crucial tool for validating the empirical privacy guarantees of differentially private algorithms. For example, it helps to identify implementation flaws in DP systems, such as incorrect noise calibration or gradient clipping mechanisms [7]. The original work [16] demonstrated how standard membership inference attacks [6] could evaluate privacy analysis algorithms, while [8] introduced worst-case poisoning examples to stress-test privacy bounds. For example, the empirical risk can be lower bounded through the membership inference performance $\epsilon \geq \ln(\text{TPR}/\text{FPR})$ [11] given ensured statistical validity.

Recent advances in DP auditing have progressed along several key dimensions: 1) **Tighter Privacy Analysis:** The work of [9] established that DP-SGD’s analysis is tight within its assumed worst-case threat model. [10] and [17] developed tighter auditing by exploiting the iterative nature of DP-SGD via auditing individual steps. Other approaches improved statistical estimation through Log-Katz confidence intervals [18] or Bayesian methods [19]. 2) **Efficiency Improvements:** [20] proposed a method to audit excluded data points without algorithm re-execution, while [21] re-uses training runs over overlapping dataset pairs to improve efficiency, though still requiring multiple runs. The recent work [11] significantly improves the auditing efficiency via the heuristic of including/excluding multiple samples in one training run and provides a theoretical bound. Recent works [13], [12] follow the similar framework further improves the tightness by leveraging f-DP definitions.

TABLE I: Notation Summary

Symbol	Description	Example/Definition
ϵ, δ	Privacy parameters	(ϵ, δ) -DP guarantee
R_g	Gradient clipping norm	$\min\{\frac{R_g}{\ g_i\ }, 1\}g_i$
z	Noise vector	$z \sim \mathcal{N}(0, I)$
σ	Noise multiplier	Scales Gaussian noise in DP-SGD
g_i	Per-example gradient	$\nabla_{\theta} l(\hat{p}(y x_i, \theta), y_i)$
D	Dataset	$\{(x_i, y_i)\}_{i=1}^n$
\mathcal{X}	Input space	Feature domain (e.g., images, text)
\mathcal{Y}	Output space	Label domain (e.g., classes, tokens)
C	Output space size	$y \in \mathcal{Y}, C := \mathcal{Y} $
B	Mini-batch	Number of samples in one batch B
θ	Model parameters	$\theta \in \mathbb{R}^d$
f_{θ}	Model function	Mapping $x \mapsto f_{\theta}(x)$
$\hat{p}(y x, \theta)$	Prediction distribution	$\text{softmax}(f_{\theta}(x))$
T	Audit trials	Number of hypothesis tests
m	Audit data size	$m > 0$
n	Training data size	$n \geq m$
α	Significance level	Type I error probability (e.g., 0.95)
ϵ_O	Optimal estimation	Assuming perfect MIA given m

Compared with these works, our work is based on the most general $O(1)$ framework [11], and improvements via tighter statistical or other DP definitions are orthogonal to ours.

B. Empirical Privacy Auditing

Not limited to auditing a theoretical lower bound for DP trained model, the general term of privacy auditing includes membership inference attack, quantifiable data proving and membership encoding. The success of $O(1)$ [11] formally proves that multi-example membership inference can be used for auditing purposes, while it also claims that it is non-trivial to leverage advanced MIA such as considering the sample hardness [2]. Similarly, the recent advanced MIA with privacy backdoor [22], [23] has been leveraged to improve $O(T)$. Quantifiable data proving [24], [25] is an instance of membership inference attack that only infers the existence of data from one owner with a certain confidence. For example, [24] creates two versions of the protected dataset with small perturbation and formulate the null hypothesis as whether the relaxed version has been used to train. Membership encoding [26], [27] is another proactive form of empirical auditing by manipulating the training process to force the model to encode membership status of its training dataset. However, these works do not directly assist DP auditing.

Compared with above empirical privacy auditing works, our work solves the problem of how to leverage sample hardness for DP auditing, and leverages insights from data proving and membership encoding for tighter theoretical DP auditing.

III. PRELIMINARIES

We will introduce background of DP training, DP auditing and our threat model for black-box $O(1)$ auditing. We summarize necessary notations in Table I.

A. Differentially-Private Training

Differential privacy (DP) provides a rigorous mathematical framework for limiting the influence of any single training example on the output of an algorithm. Formally, a randomized

algorithm is said to satisfy DP if the inclusion or exclusion of any single data point does not significantly change the distribution of its outputs. This is captured by the following definition.

Definition 1 (Differential Privacy): A randomized algorithm \mathcal{A} satisfies (ϵ, δ) -differential privacy if for any two datasets D and D' differing in at most one entry, and for all measurable subsets S of the output space of \mathcal{A} , we have

$$\Pr[\mathcal{A}(D) \in S] \leq e^{\epsilon} \Pr[\mathcal{A}(D') \in S] + \delta.$$

DP training algorithms, such as Differentially Private Stochastic Gradient Descent (DP-SGD), ensure the above definition by carefully injecting noise during the training process. Specifically, DP-SGD modifies standard SGD by clipping per-sample gradients to limit individual influence and then adding calibrated Gaussian noise to the aggregated gradients before updating the model parameters. This process guarantees that the final model's behavior does not rely heavily on any single training example, thereby satisfying (ϵ, δ) -differential privacy, which we also call it analytical privacy budget or the worst-case privacy upper bound.

For example at the iteration t , SGD optimization updates model as $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \cdot \mathbf{G}_t$, and DP-SGD privatizes the updates:

$$\text{Non-DP: } \mathbf{G}_t = \frac{1}{|B|} \sum_{i \in B} \mathbf{g}_i \quad (1)$$

$$\text{DP-SGD: } \mathbf{G}_t = \underbrace{\frac{1}{|B|} \sum_{i \in B} \min\left\{\frac{R_g}{\|g_i\|}, 1\right\} \mathbf{g}_i}_{\text{clipped gradients}} + \underbrace{\sigma R_g z}_{\text{noise scaled to clipping norm } R_g} \quad (2)$$

\mathbf{g}_i is the gradient of the loss with respect to the i -th sample in the mini-batch B , R_g is the clipping norm, σ is the noise multiplier, and $z \sim \mathcal{N}(0, I)$ is standard Gaussian noise.

In general, the DP training algorithm $\mathcal{T}(D, \theta_0 | \epsilon, \delta)$ takes the private dataset D and the initialized model θ_0 as input and outputs a trained model θ , which is claimed to satisfy (ϵ, δ) -DP guarantee that is independent on either D or θ_0 .

B. Privacy Auditing for DP Training

Differential privacy by its definition bounds the strongest adversary's capability to infer if \mathcal{T} is trained on D or D' differing by one sample. Without loss of generality, suppose guessing D' when the model is trained on D as false positive, the false positive α and false negative β are bounded by the privacy region [28]:

$$\mathcal{R}(\epsilon, \delta) = \{(\alpha, \beta) \mid \alpha + e^{\epsilon} \beta \geq 1 - \delta \wedge e^{\epsilon} \alpha + \beta \geq 1 - \delta \wedge \alpha + e^{\epsilon} \beta \leq e^{\epsilon} + \delta \wedge e^{\epsilon} \alpha + \beta \leq e^{\epsilon} + \delta\} \quad (3)$$

Ideally, ϵ can be computed by fixing δ given α and β via Equation (3). However, since the minimum possible values of α and β are hard to be computed in closed form, empirical estimates are necessary.

The essence of auditing algorithm is a hypothesis test to distinguish the model is trained on D or D' given multiple

empirical observations. In the i^{th} observation, the auditor makes a guess on the membership status with respect to a sample z_i as \hat{S}_i given the model θ trained by \mathcal{T} :

$$\hat{S}_i = f_{\text{MIA}}(I_i), \text{ where } I_i = f_{\text{SCORE}}(z_i, \theta) \quad (4)$$

A widely used scoring function is the negative loss $f_{\text{SCORE}}(z_i) = -\mathcal{L}(z_i, \theta)$, which higher score indicates lower loss and higher chance to be included in training dataset. Thus, the membership inference function $f_{\text{MIA}}(I_i)$ output binary prediction given a threshold on scores.

Given guesses and the ground truth membership S for multiple observations, it derives the confidence intervals for α and β and then convert to transfer the lower bound of ϵ given δ at a certain confidence level:

$$\epsilon_T = f_{\text{EST}}(S, \hat{S} | \delta, 1 - \beta) \quad (5)$$

For example, after obtaining the empirical lower and upper bounds using binomial proportion confidence interval for α and β . [9] uses the Clopper-Pearson method to find the upper bound of errors and derives the empirical lower bound of ϵ as:

$$\epsilon_L = \max \left\{ \ln \left(\frac{1 - \bar{\alpha} - \delta}{\beta} \right), \ln \left(\frac{1 - \bar{\beta} - \delta}{\alpha} \right), 0 \right\} \quad (6)$$

which is shown to be tight in $O(T)$ auditing for malicious dataset attack with $D_{\setminus z_i} = \emptyset$ which matches the worst-case in DP definition. Recent works improve f_{EST} tighter lower bound by using different confidence intervals [18], Bayesian techniques [19] or auditing in different privacy definitions [10].

In general, the auditing algorithm $\mathcal{A}(\mathcal{T} | \delta, 1 - \beta)$ takes the input as the implementation of training algorithm \mathcal{T} and outputs the estimated privacy risk lower bound ϵ_L under a certain confidence $1 - \beta$ given a fixed δ . We note that the output of \mathcal{A} is also influenced by the input of \mathcal{T} (the training data D and the model θ_0) and the its hyper-parameters such as η or B in Equation (2). The greatest privacy lower bound that \mathcal{A} can estimate is ϵ_O when all guesses \hat{S} are correct, which is limited by the statistical power given T observations. For example, $\epsilon_O = 5.6$ with $T = 1,000$ trials for Equation (6). For an effective DP auditing, we need $\epsilon_L \gtrsim \epsilon$.

C. $O(1)$ Framework Overview

Instead of repeating training \mathcal{T} and inference f_{MIA} for T pairs of $D \simeq D'$ for obtaining T observations, $O(1)$ framework [11] is efficient as it obtains $m \approx T$ observations in one run by excluding or including each of m audit samples. Given $n = |D|$, we have:

- The auditor samples a membership encoding vector $S \in \{-1, 1\}^m$ with each element independently drawn from the $\text{Bernoulli}(\frac{1}{2})$ distribution, with $0 \leq m \leq n$ and set $S_i = 1$ for $i \in [n] \setminus [m]^1$.
- The trainer trains $\theta = \mathcal{T}(\theta_0, D_{\text{IN}})$ with subset $D_{\text{IN}} = \{z_i | S_i = 1, z_i \in D\}_{i=1}^m$ and outputs the sorted confidence $I = \{f_{\text{SCORE}}(z_i, \theta)\}_{i=1}^m$.

¹ $[n] = \{0, 1, \dots, n-1\}$

- The auditor makes membership predictions $\hat{S} = \{f_{\text{MIA}}(I_i)\}_{i=1}^m \in \{-1, 0, 1\}^m$ with $+1(-1)$ as the positive (negative) prediction and 0 as abstention for unconfident guesses.
- The estimated risk lower bound is $\epsilon_L = f_{\text{EST}}(S, \hat{S})$.

In f_{EST} , the number of correct guesses $W := \sum_i^m \{0, \hat{S}_i \cdot S_i\}$ can be estimated over m guesses. Suppose the null hypothesis is $H_0 : \mathcal{T}$ is (ϵ, δ) -DP, the main result of O(1) in Theorem 2 indicates that if H_0 stands, $W \leq \frac{r \cdot e^\epsilon}{e^\epsilon + 1} + O(\sqrt{r})$ with high probability. Otherwise, we reject H_0 and believe \mathcal{T} violates the claimed DP. Then ϵ_L can be estimated by converting this hypothesis test into the confidence interval by finding the largest ϵ that we can reject at a desired confidence. The optimal estimation ϵ_O in O(1) scales up with the number of observation m .

Theorem 2 (Main Result of [11]): Let $(S, \hat{S}) \in \{-1, +1\}^m \times \{-1, 0, +1\}^m$ be the membership status and guess. Assume the training algorithm \mathcal{T} satisfies (ϵ, δ) -DP. Let $r := \left| \{x \in \hat{S} \mid x \neq 0\} \right|$ to be the number of guess, then for $v \in \mathbb{R}$,

$$\mathbb{P} \left[\sum_i^m \max\{0, \hat{S}_i \cdot S_i\} \geq v \right] \leq \beta + \alpha \cdot m \cdot \delta, \quad (7)$$

where

$$\beta = \mathbb{P}_{\tilde{W}^*} [\tilde{W} \geq v], \tilde{W}^* := \text{Binomial}(r, \frac{e^\epsilon}{e^\epsilon + 1}) \quad (8)$$

$$\alpha = \max_{i \in [m]} \frac{2}{i} \mathbb{P}_{\tilde{W}^*} [v > \tilde{W} \geq v - i]. \quad (9)$$

Orthogonal to the $O(1)$ framework, auditor can construct the audit dataset $D_{\text{audit}} = \{\tilde{z}_i\}_{i=1}^m \subseteq D$ with each canary crafted from $\tilde{z}_i = f_{\text{CRF}}(\tilde{z}_i)$. For example, a widely applied way is mislabeling the original sample with a random label [11].

D. Threat Model

1) $O(1)$ Auditing Game: In the auditing game, there are two parties: the *trainer* curates the DP training implementation and runs \mathcal{T} while the *auditor* performs auditing algorithm \mathcal{A} . The trainer claims that \mathcal{T} satisfies (ϵ, δ) -DP, but is potential to violate the given DP if the implementation is flawed. Therefore, the auditor aims to estimate the empirical risk lower bound ϵ_L for verifying the integrity of \mathcal{T} . Formally, we summarize the whole auditing algorithm as:

$$\mathcal{A} = f_{\text{CRF}} \circ \text{Bernoulli}(\frac{1}{2}) \circ \mathcal{T} \circ f_{\text{SCORE}} \circ f_{\text{MIA}} \circ f_{\text{EST}},$$

where the auditor performs all components except the training algorithm \mathcal{T} , and the auditor only query one run of \mathcal{T} . Multiple training runs increase the number of empirical observations, which further improve the lower bound estimation.

2) Limited Auditing Capability: Before training, the auditor will send the dataset to the trainer with a designated model architecture. Besides, we assume a strict and practical capabilities for auditors as follows:

- **No manipulation on implementation:** the auditor does not interfere the DP training implementation, including gradient calculation or hyper-parameters. Existing works

may modify the normal DP training step for maximizing auditing capability, such as by ignoring updates of non-audit samples or crafting the real gradients, while we do not assume such capability.

- **Normal training configuration:** we assume all configurations are normal without requirement of specific model initialization or hyper-parameter tuning (including learning rate or batch size). While existing works shown that using privacy-backdoor [22], [23], an extremely large learning rate or sub-sampling ratio helps to improve ϵ_L , we do not force such assumption for a practical use.
- **Black-box access:** the auditor cannot obtain the target model's parameters and only audit through a black-box access. Existing works show that black-box is much more challenging than white-box in O(1) framework, and can only achieve $\epsilon_L \approx 1.3$ [11], [29] for $\epsilon = 4$ given confidence $\alpha = 0.95$ and $m = 1,000$.
- **Final model access:** the auditor can only get access to the final model in our setting, while existing works [9] assume auditors can access all intermediate models.

IV. NO FREE LUNCH FOR O(1) AUDITING

Even though O(1) framework in Section III-C is efficient, it sacrifices the auditing power due to two main challenges. For illustrating the key challenges, we train CNN model on CIFAR10 dataset from scratch and perform a grid-search over hyper-parameters that may influence both auditing power and model utility as shown in Figure 1.

We consider two versions of canary function f_{CRF} in O(1): 1) *in-distribution* uses the original CIFAR10 subset as D_{audit} ; 2) *misabeled* flips the true labels, and one variant with stronger f_{MIA} 3) *poisoned* applies privacy backdoor [22], [23].

A. Inevitable Data Dependency

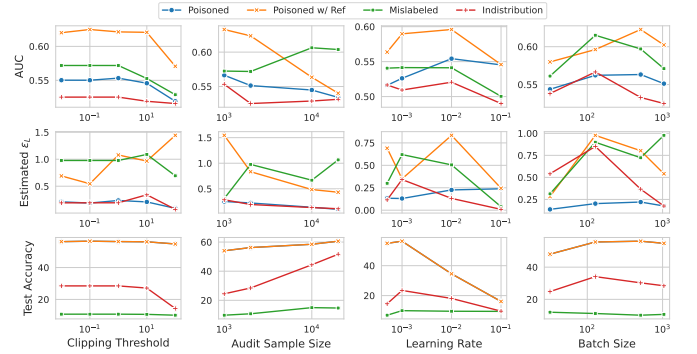
While in O(T) auditing, each observation is obtained from *independent* training runs. Observations obtained in O(1) are partially *dependent*, as including or excluding one sample influences other audit samples. And such data dependency influences in both training \mathcal{T} and inferences $f_{\text{SCORE}} \circ f_{\text{MIA}}$. For example, similar samples with the same labels in D_{audit} are essentially correlated by the core feature mapping to the label. Thus, these samples have a regularization effect [30] to each other, reducing the sample-specific memorization in \mathcal{T} . Consequently in the inference phase $f_{\text{SCORE}} \circ f_{\text{MIA}}$, member scores are less distinguishable, resulting in weak membership inference attack (MIA).

As the result shown in the second column of Figure 1(a), both the estimated ϵ_L and MIA AUC decrease with larger m when audit canaries are real samples with correlation. For mislabeled baseline with random labels (less dependency), ϵ_L scale up with m while also fluctuates. In summary, the inevitable data dependency limits auditing power of O(1) framework from two aspects:

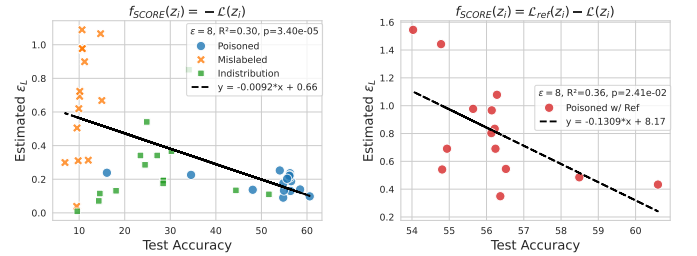
- **Limited number of observation:** the optimal estimation ϵ_O (when f_{MIA} is 100% accurate) is **not** positively correlated

with m as expected; thus, it is **impossible** to spot privacy violation when $\epsilon_O < \epsilon$, no matter how strong is f_{MIA} .

- **Limited MIA capability:** the dependent scores impair the MIA capability which directly result in weaker ϵ_L .



(a) ϵ_L fluctuates with hyper-parameters, not scaling with m .



(b) Trade-off between Utility and Auditing for O(1)

Fig. 1: Preliminary analysis for O(1) Auditing. Fig (b) is summarized from checkpoints in Fig (a). By default, we use $m = 2,000, B = 2,000, \eta = 1e-3, E = 100$ and the theoretical epsilon is $\epsilon = 8$. For *Poisoned* baseline, the adversarial model initialization is warmed up on an extra auxiliary dataset with 5,000 samples and 50 epochs.

B. Implicit Conflict between Auditing and Utility

We are the first to reveal that there is an implicit conflict between auditing and utility, which also limits the auditing power of O(1). As shown in Figure 1(a), ϵ_L is sensitive to hyper-parameter choices. While existing works mainly focus on tighter estimation and opt for a larger learning rate η or batch size B to maximize the auditing result, once the auditor attempts to balance the auditing-utility trade-off [29], auditing power is clearly compensated as the negative trend shown in Figure 1(b). Specifically, *poisoned* seems has higher test accuracy and higher ϵ_L with advanced $f_{\text{SCORE}} \circ f_{\text{MIA}}$, it requires extra **non-private** auxiliary dataset, which is not comparable to *mislabeled* and *in-distribute*. Without such assumption, *mislabeled* has **no meaningful** utility as shown in Figure 1(b)-left while achieves higher ϵ_L than *in-distribute*.

C. Our Key Insights from the Generalization Perspective

We observe that the key to addressing the above two challenges lies in the boundary between sample-specific memorization and generalization.

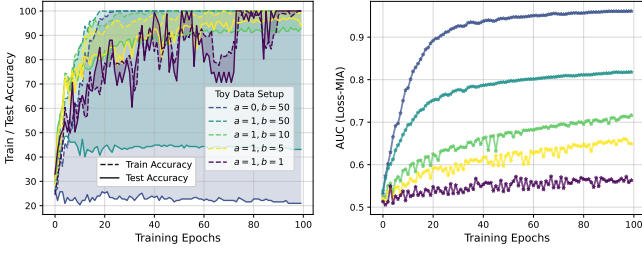


Fig. 2: Interplay between sample-specific memorization (less data correlation for auditing) and generalization (more data correlation for utility).

Setup for Toy Experiments. To investigate the interplay, we construct a toy dataset that serves as an abstracted instance of real-world data, without loss of generality. In our setup, each training sample ideally consists of informative features that determine its correct label. Given a label y , we generate the corresponding sample as

$$\mathcal{X}_y = a\mathcal{N}(y, \sigma_0^2) + b\mathcal{N}(0, \sigma_0^2) \in \mathbb{R}^d, \quad (10)$$

where the first term represents the core feature that depends solely on the label y , and the second term adds sample-specific Gaussian noise with the basic noise scale $\sigma_0 = 0.1$. We use $a \in \{0, 1\}$ to toggle the presence of label-dependent features, and $\alpha = 1$ simulates real-data with correlated features and labels. The second term $b \geq 0$ controls the level of sample-specific variation. We optimize two-layer neural network with ReLU activations with SGD on the toy data for 100 epochs.

Key Insight I. Uncorrelated Data for Less Dependency. Fixing the noise level $b = 50$, Figure 2 shows that $a = 0$ with correlated labels leads to larger generalization gap (left) and higher MIA AUC (right) than $a = 1$ which is more like real-world data. Fixing the correlation level $a = 1$, larger relative sample-specific noise level b/a enhances MIA. This reveals the fact that auditing power depends on the sample-specific noise and the correlation level between feature and label impairs the relative sample-specific signal. Therefore, we are motivated to design f_{CRF} by constructing D_{audit} with uncorrelated data for reducing dependency in $\mathcal{T} \circ f_{\text{SCORE}} \circ f_{\text{MIA}}$.

Key Insight II. Decoupling Auditing from Utility. Real-data definitely has correlated feature and labels for meaningful test performance. Given $a = 1$ in Figure 2, the generalization gap decreases and the testing accuracy gets better with smaller sample-specific noise level b , while the MIA AUC gets worse. The implicit conflict is clear, as the auditing power relies on the sample-specific noise, which hurts the model generalization capability for good utility. This tension motivates us to challenge the necessity of coupling utility and auditing in $O(1)$. Our intuition is two fold:

- *Decoupled Criterion.* As mentioned in Section III-A, the definition of DP rely on nothing but \mathcal{T} . Therefore, as we summarized in Table II, meaningful model utility for auditors is not a necessary criterion for **Case I** where the goal is to diagnose DP implementation with

TABLE II: Overview of $O(1)$ auditing variants without manipulation training algorithm and black-box access to the final model. (● indicates ‘Yes’, ○ indicates ‘No’, ∅ means $D_{\text{audit}} = D$ and $n = m$)

O(1) Audit Variants	Data Correlation		Need $D_{\text{aux}}?$	Criterion	
	D_{audit}	$D \setminus D_{\text{audit}}$		Utility	Auditing
In-distribution	●	●	○	●	○
Mislabeled Data	○	●	○	●	●
Poisoned Model	●	●	●	●	●
UniAud (Data-Independent)	○	∅	○	○	●
UniAud++ (Data-Dependent)	○	●	○	●	●

TABLE III: Canaries in DP auditing (● indicates ‘Yes’)

Canary	Real-data?	Need $D_{\text{aux}}?$	Optimization?	For $O(1)?$	For black-box?
Blank	○	○	○	○	●
In-distribution	●	○	○	●	●
Mislabel	●	○	○	●	●
PGD	●	○	●	○	○
Orthogonal	●	●	●	○	○
ClipBKD	●	●	●	○	○
Ours	○	○	○	●	●

$\epsilon_L = \mathcal{A}(\mathcal{T}|\ast)$. The auditing-utility trade-off should be considered as criterion in **Case II** where $\epsilon_L = \mathcal{A}(\mathcal{T}|D, \theta)$ is expected to be conditioned on a specific run.

- *Decoupled Objective.* For **Case II** with coupled criterion, our idea is to separate training objectives for utility and training for auditing. The utility is optimized with main task objective, while we formulate auditing goal as a membership encoding problem.

Based on our key insights, we will improve the $O(1)$ framework with a data-independent framework UniAud for Case I in Section V, and a data-dependent $O(1)$ framework UniAud++ for Case II in Section VI.

V. DATA-INDEPENDENT $O(1)$ AUDITING FRAMEWORK

In this section, we propose data-independent framework UniAud, which aims to improve auditing power within $O(1)$ complexity by reducing data dependency via carefully designed f_{CNR} , f_{SCORE} and f_{MIA} .

A. Building Canaries for Membership Encoding

We first resolve data dependency problem with better f_{CNR} .

1) *Rules of Thumb for Audit Canaries:* As summarized in Table III, many canary strategies used in $O(T)$ auditing [8], [10] are not applicable to $O(1)$. For designing the synthetic audit data, here are rules of Thumb for canaries in $O(1)$ framework: **a) Independence:** Each canary should contribute independently to the audit, unlike methods such as PGD [10], Orthogonal [10], or ClipBKD [8], which assume $O(n)$ settings. **b) Easy-to-Spot:** Canaries should be easily distinguishable to enable tighter auditing bounds. **c) Scalable Generation:** $O(1)$ auditing benefits from many canaries in a single run, making per-sample optimization or reliance on auxiliary data impractical.

2) *Uncorrelated Pair-Matching as Membership Encoding:* Based on Section IV-C, the essence of auditing is better membership encoding, with the goal to minimize loss of member

samples while maximize the loss of non-member samples. Unlike prior work in non-DP settings, our auditor cannot modify the training procedure; thus the key is constructing the joint space of audit data for better membership encoding.

We propose to use **uncorrelated** pair-matching task for membership encoding, with the objective to minimize the probability of matching a random feature to a random label:

$$\mathcal{L}_{\text{ME}} = \frac{1}{m} \sum_{i=1}^m \mathbb{I}(S_i = 1) \mathcal{L}_{\text{CE}}(f_{\theta}(x_i), y_i), \quad (11)$$

where x_i and y_i are drawn independent from \mathcal{X} and \mathcal{Y} . Without loss of generality, we instantiate as a multi-label classification problem with $\mathcal{X} \subset \mathbb{R}^d$ and $\mathcal{Y} \in [C]$ of C classes, leaving auto-regressive extensions in Section VII.

We leverage this high-dimensional sparsity by independently sampling m audit points uniformly from all possible input-label combinations from $\mathcal{X} \times \mathcal{Y}$, indicating that each $(x_i, y_i) \in D_{\text{audit}}$ is very likely separated from others. As $|\mathcal{X}|$ and $|\mathcal{Y}|$ grow, the joint space expands exponentially, making it virtually impossible for the model to generalize as for unseen non-member samples $(x_i, y_i) \in D_{\text{audit}}$ with $S_i = -1$:

$$\mathbb{E}_{(x,y) \sim \mathcal{X} \times \mathcal{Y}} [p_{\theta}(y|x)] \approx \frac{1}{C} \implies \mathcal{L}_{\text{CE}}(x_i, y_i) \approx \log C$$

creating an inherent loss gap between members ($\mathcal{L}_{\text{CE}} \rightarrow 0$) and non-members ($\mathcal{L}_{\text{CE}} \approx \log C$).

3) *Synthetic Canary Generation*: Now we discuss concrete construction for \mathcal{X} and \mathcal{Y} . We notice that pure random \mathcal{X} with uncorrelated labels \mathcal{Y} is the ideal solution that satisfies above three rules and fully expand the joint space. Especially when \mathcal{X} is pure random (as $a = 0$ in Figure 2), samples are more separated than natural images which center around true labels, leading to faster convergence rate [31] than only random label. We propose synthetic canary generation with two variants in Section V-A3. Specifically for orthogonal features, output features are approximately orthogonal when $n > d$.

Algorithm 1 Synthetic Data Generation for $D_{\text{audit}} \in \mathcal{X} \times \mathcal{Y}$

Require: Number of audit samples m ; \mathcal{X} with mode $\in \{\text{gaussian, orthogonal}\}$, feature dimensions d , and sample-noise scalar σ_0 ; $\mathcal{Y} = \text{Uniform}([C])$ with C classes

Ensure: Synthetic dataset $D_{\text{audit}} = \{(x_i, y_i)\}_{i=1}^m$

```

1: if mode == orthogonal then
2:   // Generate orthonormal matrix  $Q$ 
3:    $Q, R \leftarrow \text{QR}(\mathcal{N}(0, 1)^{d \times d})$ 
4:   // Normalize coefficient vectors
5:    $U \leftarrow \mathcal{N}(0, 1)^{m \times d}$ 
6:    $U \leftarrow U / \|U\|$  row-wise
7:    $\{x_i\}_{i=1}^m \leftarrow U \cdot Q^{\top}$ 
8: else if mode == gaussian then
9:    $\{x_i\}_{i=1}^m \leftarrow \mathcal{N}(0, \sigma_0^2)^{n \times d} \cdot \sigma_0$ 
10:  $\{y_i\}_{i=1}^m = \{y_i | y_i \in \mathcal{Y}, i \in [m]\}$ 
11: return  $D_{\text{audit}}$ 
```

4) *Model Architecture*: By default for classification task, we use a 2-layer neural network with ReLU activations. Following previous work [31], a two-layer ReLU neural network with $2n + d$ parameters can universally represent any function defined on a sample of size n in d -dimensional space. In general, the synthetic D_{audit} is compatible with different model architectures with input dimension d and output size C .

B. Inclusive auditing by self-comparison

Algorithm 2 Data-Independent O(1) Auditing (synthetic canary + self comparison)

- 1: **Input:** training algorithm \mathcal{T} , canary distribution $\mathcal{X} \times \mathcal{Y}$, audit confidence level $1 - \beta = 0.95$
 - 2: **Output:** estimated privacy lower bound ϵ_L
 - 3: Auditor gets audit set $D_{\text{audit}} \leftarrow f_{\text{CRN}}(m, \mathcal{X} \times \mathcal{Y})$ and comparison set as $D_{\text{comp}} \leftarrow \{x_i, y' | y' \in \mathcal{Y}\}_{i=1}^m$
 - 4: Trainer trains on audit data $\theta = \mathcal{T}(D_{\text{audit}} | \theta_0)$
 - 5: Sample m evaluation samples D_{mia} , with membership status $S_i \in \{-1, +1\}$ drawn from $\text{Bernoulli}(\frac{1}{2})$:
 - 6: $z_i \leftarrow \begin{cases} D_{\text{comp}}[i]; \tilde{z}_i = D_{\text{audit}}[i] & S_i = -1 \\ D_{\text{audit}}[i]; \tilde{z}_i = D_{\text{comp}}[i] & S_i = 1 \end{cases}$
 - 7: Compute scores for D_{mia} as $I = \{f_{\text{SCORE}}(z_i, \theta)\}_{i=1}^m$ where $f_{\text{SCORE}}(z_i) = \mathcal{L}(f_{\theta}, \tilde{z}_i) - \mathcal{L}(f_{\theta}, z_i)$
 - 8: f_{MIA} sorts I and make r guesses with $r/2$ samples with highest (w.r.t. lowest) I_i as $\hat{S} = 1$ (w.r.t. $\hat{S} = -1$)
 - 9: Estimate privacy risk lower bound $\epsilon_L = f_{\text{EST}}(S, \hat{S}, \delta, 1 - \beta)$
-

Uncorrelated synthetic canary constructed to enlarge the generalization gap with the pair-matching task makes scores less dependent and more distinguishable. Now we further improve MIA capability with better $f_{\text{SCORE}} \circ f_{\text{MIA}}$. The key intuition that we can further improve MIA for auditing is attributed to the label independency Property 1.

Property 1: Label Independency: For all $x \in \mathcal{X}$ and $y, y' \in \mathcal{Y}$, the function f_{CRN} in Section V-A3 satisfies

$$\Pr_{\mathcal{X} \times \mathcal{Y}} [y | x] = \Pr_{\mathcal{X} \times \mathcal{Y}} [y' | x].$$

Existing works [11] use $f_{\text{SCORE}} = -\mathcal{L}(f_{\theta}(x_i), y_i)$ and then sort scores for f_{MIA} given negative and positive thresholds, resulting in guesses dependent on sample-hardness of D_{audit} . This issue is unique for O(1) because each observation in O(T) is only obtained by differing one sample. And as the original work points out, it is non-trivial to leverage advanced MIAs with reference models to calibrate each sample hardness.

Instead, we propose Algorithm 2 which makes f_{MIA} independent among audit samples by calibrating the score f_{SCORE} with a counterfactual calibration. Specifically, we make the following modifications compared to the original O(1) framework: a) We train on the full set of D_{audit} and construct non-member set by independently sample a fresh label from \mathcal{Y} . b) We independently flip the observation target via $\text{Bernoulli}(\frac{1}{2})$, ensuring an unbiased f_{MIA} . c) We only need one-side threshold because the counterfactual calibrated f_{SCORE} is symmetric. Thus, we have Theorem 3 with the key prerequisite of Property 1, with proof in Appendix.

Theorem 3: Assume training algorithm \mathcal{T} satisfies (ϵ, δ) -DP, and f_{CNR} in Section V-A3 is used in auditing Algorithm 2, then the inequality in Theorem 2 hold for $(S, \hat{S}) \in \{-1, +1\}^m \times \{-1, 0, +1\}^m$ in Algorithm 2.

VI. DATA-DEPENDENT O(1) AUDITING FRAMEWORK

In this section, we propose data-dependent framework UniAud++ with aims to improve the auditing-utility trade-off given good efficiency with O(1) complexity for Case II. This scenario is practical when auditor aims to estimate the conditioned empirical risk for a specific DP training run given data D and model θ , or when one DP training is expensive thus θ is required to have meaningful utility.

A. Multi-Task for Decoupled Objective

While the conflict between utility and auditing is inherent as shown in Section IV-C, we note that it is because the in typical auditing training, the two training objectives are disentangled. Specifically, m original or mislabeled canaries still share the same target space \mathcal{Y} with the rest $n - m$ non-audit samples in the training set D .

Our key insight is to decouple the objective of auditing and utility during DP training by augmenting training dataset with an extra feature and label space. Thus, we transfer the data distribution as $\mathcal{X} \times \mathcal{G} \rightarrow \mathcal{Y} \times \mathcal{E}$, where \mathcal{G} and \mathcal{E} denote the trigger space and tag space for membership encoding. Thus, combining the membership encoding objective in Equation (11) with the typical main task, we have the training objective as

$$\mathcal{L} = \underbrace{\mathcal{L}_{\text{CE}}(f_{\theta}(x_i), y_i)}_{\mathcal{L}_{\text{main}}} + \lambda \underbrace{\mathcal{L}_{\text{CE}}(f_{\phi}(x_i), e_i) \mathbb{I}(S_i = 1)}_{\mathcal{L}_{\text{ME}}}, \quad (12)$$

where λ is the coefficient for balancing encoding strength. Here, ϕ denotes an additional linear head on top of the shared encoder, following the standard multi-task setup with separate heads for distinct objectives. This allows joint optimization of the main task and membership encoding.

For optimizing objective, we construct dataset in Section VI-A, which returns D_{multi} for calculating Equation (12). The self-comparison framework Algorithm 2 is seamlessly integrated with the returned D_{audit} with D_{comp} by replacing notation of $\mathcal{X} \times \mathcal{Y}$ with $\mathcal{G} \times \mathcal{E}$ and f_{SCORE} is calculated on the separated encoding loss \mathcal{L}_{ME} .

B. Multi-Bit Membership Encoding

Essentially, the joint space $\mathcal{G} \times \mathcal{E}$ induces sample-specific patterns through unique $\langle g_i, e_i \rangle$ pairs, enabling precise membership encoding. If two audit samples have same trigger-tag pair, they are not independent enough. According to the birthday attack [32], the collision probability that two samples have same $\langle g_i, e_i \rangle$ is $C(|\mathcal{E}|, m) \approx \frac{m^2}{2|\mathcal{E}|}$ for $1 \leq m \leq \sqrt{2|\mathcal{E}|}$. Thus, the ideal number of audit samples for independent membership encoding is $m \ll \sqrt{2|\mathcal{E}|}$.

For a linear projection head mapping d_h -dimensional features to $|\mathcal{E}|$ classes, this requires $d_h \times |\mathcal{E}|$ parameters—a potential bottleneck for utility and efficiency of DP training

Algorithm 3 Multi-Task Dataset Construction for UniAud++

Require: Real dataset $D = \{(x_i, y_i)\}_{i=1}^n$; number of audit samples m ; trigger generator \mathcal{G} ; tag generator $\mathcal{E} = \text{Uniform}([C_e])$

Ensure: D_{multi} for training, $D_{\text{audit}}, D_{\text{comp}}$ for auditing

- 1: Sample m indices $\mathcal{I}_{\text{audit}} \subset [n]$ uniformly at random
- 2: Initiate empty dataset $D_{\text{multi}}, D_{\text{audit}}, D_{\text{comp}}$
- 3: **for** $i = 1$ to n **do**
- 4: **if** $i \in \mathcal{I}_{\text{audit}}$ **then**
- 5: Generate trigger $g_i \sim \mathcal{G}$ and sample tag $e_i \sim \mathcal{E}$
- 6: Generate a fresh tag $e'_i \sim \mathcal{E}$
- 7: Apply trigger $x_i \leftarrow x_i \circ g_i$ and set $S_i \leftarrow 1$
- 8: Add (x_i, y_i, e_i, S_i) to D_{audit}
- 9: Add (x_i, y_i, e'_i, S_i) to D_{comp}
- 10: **else**
- 11: $e_i \leftarrow \text{None}, S_i \leftarrow 0$
- 12: Add (x_i, y_i, e_i, S_i) to D_{multi}

at scale. We propose multi-bit encoding which encodes each trigger with multiple tags $\{e_i \in \mathcal{E}\}_{i=1}^H$, thus expanding the number of unique tags to by a factor of $\binom{\mathcal{E}}{H}$. Suppose $|\mathcal{G}| \cdot |\mathcal{E}|$ trigger tag combinations can cover m sample-specific encoding, it requires $|\mathcal{E}| \cdot d_h$ number of weights in its linear head. We can reduce the number of weights required for m samples by a factor of $1/\binom{\mathcal{E}}{H}$, which achieves its minimum when $H = |\mathcal{E}|/2$ given even $|\mathcal{E}|$ is even; and $(|\mathcal{E}| + 1)/2$ and $(|\mathcal{E} - 1|)/2$ for odd $|\mathcal{E}|$.

VII. EXTENDING TO LANGUAGE MODELING

Our framework UniAud and UniAud++ are general and scalable to other tasks besides image classification. Now we exemplify that both proposed frameworks can be seamlessly extended to language models (LM).

For LMs in UniAud, we adapt the trigger-tag pair as a prefix-suffix pair matching problem. Formally, in the auto-regressive task of LM, we denote the text sequence as $\mathbf{x} \in \mathcal{X}$ with prefix $\mathbf{x}_p \circ \mathbf{x}_s$. We propose to leverage LM's capability on learning the repetition pattern to maximize the membership encoding, which is different from a recent work which uses random tokens [29]. To avoid influence of existing tokens as [33], we augment the token vocabulary from V with V_{new} . For each $\mathbf{x}_i \in D_{\text{audit}}$ for $i \in [m]$, we independently sample only two random tokens from the new token set V_{new} , and repeat them with $|\mathbf{x}_p|$ and $|\mathbf{x}_s|$ times to create one canary. A freshly sampled token is sampled from the same V_{new} for $\mathbf{x} \in D_{\text{comp}}$ by fixing the t_p -length prefix.

Then we can easily generalize above solution by adding $n - m > 0$ samples of normal text with $S_i = 0$ in UniAud++ for meaningful utility. The role of token embeddings V_{new} in our LM variant is equal to the extra task head for optimizing the second term in Equation (12) for classification. One difference is that there is no need to merge such pair pattern in $n - m$ normal samples because the optimization over new tokens in V_{new} and normal text tokens in V are naturally separated.

VIII. EXPERIMENTS

As the proposed framework is unified to different inputs, we evaluate UniAud in **Case I** and UniAud++ in **Case II** across image classification tasks and language modeling tasks.

A. Experimental Setups

Tasks and Dataset. In image classification task, for covering general types of model architectures, we evaluate convolutional neural network [34] and transformer-based model of ViT², both of which are widely implemented in standard DP training benchmarks [34], [35], [36], [33], [37]. To investigate the influence of model initialization, we consider both training from scratch, and fully fine-tuning over pre-trained parameters. By default, the ViT encoder [38] is pre-trained on ImageNet-21k, thus it is capable to capture general vision representation for classification. Besides our synthetic audit data, we evaluate on standard CIFAR110 dataset used in existing auditing benchmark and additionally include the more difficult CIFAR100. We also evaluate on GTSRB, which contains real-world traffic sign images with high inter-class similarity and varying image quality.

In language modeling (LM) task, we follow previous work [29] and use GPT-2 family as an representative model architecture. It should be noted that the improvement of auditing framework does not depend on specific architecture. And following [29], we evaluate on PersonaChat dataset, and additionally include a subset of PubMed data in 2023, after GPT-2 pre-training data’s cut off date in 2019, both datasets are widely used in previous DP training works [39], [40].

Baseline. In general, the key difference between data-independent (Case I) and data-dependent (Case II) setting is that $m = n$ in Case I for less influence from non-audit data for stronger auditing, while $m < n$ in Case II for maintaining acceptable utility. For both cases, we compare with the following black-box O(1) variants with either different canary types and advanced MIA components.

- 1) *In-distribution* is the standard baseline of O(1) framework [11] under black-box setting with a random subset of samples in the original training dataset as audit data.
- 2) *Mislabel* represents a popular canary construction as also used in O(1) [11] which keeps the original features but randomly flips the labels of the m audit samples. As for LM tasks, it equals to a recent work [29] that constructs canary with in-distribution prefix and random tokens as suffix.
- 3) *Poisoned* is adapted from advanced MIA [22], [23] by crafting the model initialization for amplifying the privacy risk of the training dataset, which represents a stronger auditing capability as it leverages extra assumptions beyond the auditing capability as we defined in Section III-D. We build the poisoned model by warming up the model on an in-distribution auxiliary subset.

- 4) *NewToken* [15] is a state-of-the-art canary construction method specifically for LMs, by using new tokens as suffix. We follow the best performed setting [15] by using random prefix and supervised fine-tuning (SFT) loss objective for DP auditing.

Ours indicates UniAud in case I and UniAud++ in case II. For case I, we by default use a 2-layer MLP with ReLU network [31] for image classification and use GPT-2 for language modeling task, because utility is not a necessary criterion. By default, we set the input feature dimension as $d_x = 10^3$, hidden states dimension as $d_h = 10^5$ and output space size $C = 10^3$ for maximizing the joint space size.

For image tasks in case II, UniAud++ is built on UniAud by decoupling the training objective and revising the way of integrating sample-specific features in training data for maintaining utility. While for LM, the training objective and way of inserting canary feature keeps the same as we find the utility degradation is low especially for $n > m$ in case II due to the naturally separate space of normal and new tokens.

By default, we report results with $\delta = 1 \times 10^{-5}$ and confidence 95%. We use the same training hyper-parameters for one type of model architectures, and we follow previous work to set the sub-sampling ratio for all DP training as 0.1.

B. Evaluation of Data-Independent Auditing

We first evaluate UniAud for data-independent auditing in **Case I** where the criterion of auditing algorithm is the tighter ϵ_L compared to the upper bound ϵ .

1) *Auditing Improvement across Privacy Budget:* We demonstrate the estimated privacy risk lower bound ϵ_L across different analytical privacy budget $\epsilon = \{1, 2, 4, 8, \infty\}$ in Figure 3. The non-DP training with $\epsilon = \infty$ in the figure acts as an indicator of the capability of f_{MIA} in the auditing algorithm \mathcal{A} . Given the optimal estimated risk ϵ_O calculated when we assume the MIA component f_{MIA} has 100% inference accuracy, thus if an auditing algorithm has $\epsilon_L < \epsilon_O$ for $\epsilon = \infty$, it indicates that the MIA f_{MIA} is weak and impairs the auditing tightness. We use orthogonal canaries with $m = 2,000$ for ours, and optimize $m = \{500, 2,000, 5,000\}$ for each baselines of real-data canaries.

First of all, we can observe that the tightest ϵ_L (shown in gray shadow) across different ϵ is achieved with our Section V-A3 and the default architecture of 2-layer MLP. The auditing performance on GPT-2 also approaches to the tightest results, while the gap between ours and all variants of $\{\text{In-distribution}, \text{Mislabel}, \text{Poisoned}\} \times \{\text{GTSRB}, \text{CIFAR10}, \text{CIFAR100}\}$ is significant. The improvement comes from two aspects: the uncorrelated canary construction f_{CNR} in Section V-A3 and the MIA capability of f_{MIA} in Algorithm 2.

We note that baselines for image tasks have larger gap between ϵ_L and optimal estimates ϵ_O when $\epsilon = \infty$, which reflects the influence of data dependency on $f_{\text{SCORE}} \circ f_{\text{MIA}}$. While Mislabel baseline in LM task has higher ϵ_L than ours, the extra randomness introduced in DP training algorithm significantly limits its power when $\epsilon \neq \infty$.

²https://github.com/huggingface/pytorch-image-models/blob/main/timm/models/vision_transformer.py

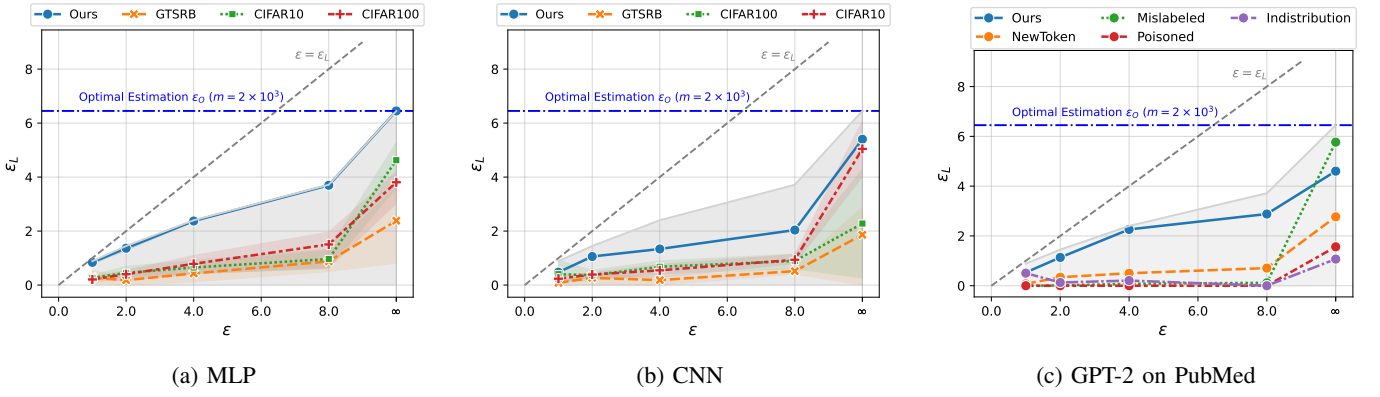


Fig. 3: Auditing improvement in Case I of data-independent setting with $\delta = 1e-5$ and 95% confidence. Gray shadows indicate the tightest ϵ_L achieved by the default MLP with orthogonal synthetic canary in UniAud across models, which outperforms the state-of-the-art black-box performance [11], [29] (e.g., $\epsilon_L \approx 1.3$ for $\epsilon = 4$), and approaches to white-box auditing performance (e.g., $\epsilon_L \approx 2.2$ for $\epsilon = 4$) but saves at least $T = 10^3$ training runs. In (a)(b), the range covered by each dataset (GTSRB, CIFAR10, and CIFAR100) summarizes ϵ_L across baselines (In-distribution, Misabeled, and Poisoned).

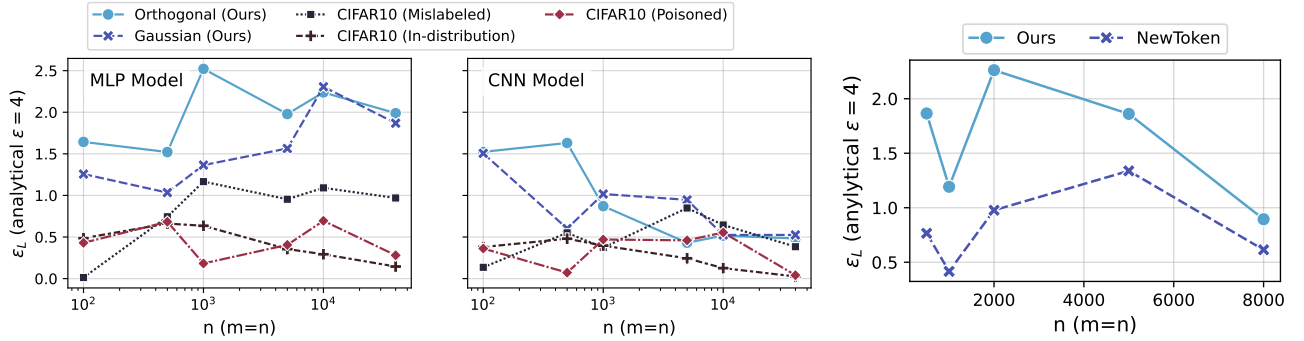


Fig. 4: Influence of m for image classification (left and middle) and language modeling (right)

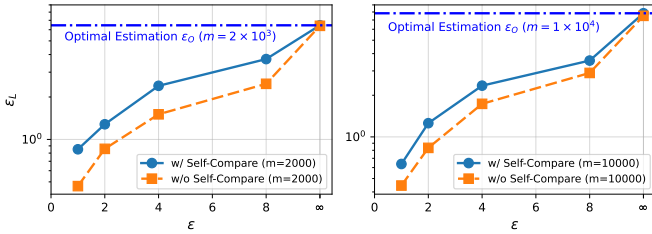


Fig. 5: Effectiveness of self-comparison with MLP model

Additionally, we set the same training hyper-parameters for all baselines and our auditing for a fair comparison. Thus, our improvement does not depend on model architecture, even lower than MLP model and GPT-2 model, ϵ_L for CNN still outperforms other baselines with real-data as canary and sorting for MIA. Specifically for each baseline in each dataset and each model combination, we demonstrate a detailed result in Table IV. We can observe that

2) *Auditing Improvement across Different m* : Now we compare UniAud with baselines across m , the number of auditing data. The size of audit data is vital to the auditing

performance, as it not only controls the statistical power as reflected in ϵ_O , and also influences the dependency level for the m observations. The larger m results in larger ϵ_O but may also involves more dependency that hurt the auditing power.

As shown in Figure 4, the two types of synthetic canary that we propose in Section V-A3 achieve a better auditing result than all baseline, as a result of uncorrelated audit data and the self-comparison framework. By comparing different baselines, we find that *Misabeled* is superior than other baselines, due to its less dependency on label space, but still inferior to our methods with both noisy features and noisy labels because noisy features are more separated in the space [31].

In the right sub-figure for language modeling task, we compare with the strongest baseline *NewToken*. For a fair comparison, we insert the same number of new tokens, keep all training hyper-parameters the same, e.g., using the same supervised fine-tuning loss by masking the loss calculation in the prefix. The key difference between ours and *NewToken* is that they use random prefix and suffix while ours use repeated tokens as prefix and suffix to construct the canary. Our implementation has achieved the their reported optimal $\epsilon_L \approx 1.3$ given $\epsilon = 4$, $\delta = 1e-5$ given 95% confidence. And

TABLE IV: Estimated privacy lower bound ϵ_L under confidence $1 - \beta = 0.95$. The column with $\epsilon = \infty$ denotes the MIA capability limit in auditing. The optimal estimation ϵ_O is derived by assuming all predictions are correct, indicating the limit of statistic power limit in auditing.

Dataset	Model	$m = 2 \times 10^3$ ($\epsilon_O = 6.45$)			$m = 10 \times 10^3$ ($\epsilon_O = 7.83$)		
		$\epsilon = 1$	$\epsilon = 8$	$\epsilon = \infty$	$\epsilon = 1$	$\epsilon = 8$	$\epsilon = \infty$
CIFAR10	CNN	0.132	0.193	1.631	0.090	0.134	0.361
	ViT-Small	0.173	0.099	0.948	0.000	1.499	0.510
	ViT-Base	0.067	0.153	0.016	0.010	0.029	0.327
CIFAR100	CNN	0.183	0.165	4.507	0.230	0.118	6.601
	ViT-Small	0.322	0.359	4.701	0.031	0.477	6.242
	ViT-Base	0.265	0.656	1.263	0.068	0.457	1.755
GTSRB	CNN	0.004	0.131	1.986	0.003	0.105	0.368
	ViT-Small	0.101	0.387	0.479	0.091	0.050	3.534
	ViT-Base	0.376	0.981	1.189	0.053	0.061	0.259
Mislabelled		$\epsilon = 1$	$\epsilon = 8$	$\epsilon = \infty$	$\epsilon = 1$	$\epsilon = 8$	$\epsilon = \infty$
CIFAR10	CNN	0.739	1.219	1.730	0.068	0.465	1.168
	ViT-Small	0.036	0.099	1.703	0.151	0.655	1.793
	ViT-Base	0.001	0.739	0.196	0.110	1.195	0.928
CIFAR100	CNN	0.038	0.291	0.286	0.068	0.052	0.047
	ViT-Small	0.228	0.708	2.881	0.084	0.641	1.722
	ViT-Base	0.310	0.745	2.571	0.059	0.884	1.215
GTSRB	CNN	0.000	0.207	0.021	0.000	0.107	0.000
	ViT-Small	0.130	0.316	0.851	0.108	0.521	4.015
	ViT-Base	0.024	0.749	0.409	0.492	0.630	0.856
Poisoned		$\epsilon = 1$	$\epsilon = 8$	$\epsilon = \infty$	$\epsilon = 1$	$\epsilon = 8$	$\epsilon = \infty$
CIFAR10	CNN	0.052	0.834	0.421	0.002	0.466	0.554
	ViT-Small	0.341	0.483	0.488	0.466	1.499	4.963
	ViT-Base	0.981	0.411	1.644	0.047	0.678	3.047
CIFAR100	CNN	0.010	1.189	5.871	0.128	1.145	6.571
	ViT-Small	0.003	1.382	4.932	0.466	1.745	5.928
	ViT-Base	0.388	1.741	5.009	0.243	0.813	5.926
GTSRB	CNN	0.523	1.124	3.506	0.033	0.958	0.976
	ViT-Small	0.000	0.153	2.338	0.000	0.993	2.655
	ViT-Base	0.000	0.158	2.414	0.580	0.595	0.958
Ours		$\epsilon = 1$	$\epsilon = 8$	$\epsilon = \infty$	$\epsilon = 1$	$\epsilon = 8$	$\epsilon = \infty$
Gaussian w/o comp	2-Layer ReLU	0.322	2.362	6.395	0.207	2.096	7.816
	CNN	0.415	1.281	4.931	0.120	1.036	6.259
	ViT-Small	0.124	1.116	6.395	0.378	1.318	7.503
	ViT-Base	0.238	1.038	6.395	0.244	1.422	7.503
Gaussian w/ comp	2-Layer ReLU	0.620	2.707	6.449	0.771	3.780	7.829
	CNN	0.579	1.281	5.395	0.171	1.423	6.816
	ViT-Small	0.413	1.306	6.449	0.173	1.446	7.834
	ViT-Base	0.174	1.365	6.449	0.521	1.601	7.834
Orthogonal w/o comp	2-Layer ReLU	0.447	3.032	6.395	0.501	3.191	7.503
	2-Layer ReLU	1.089	3.059	6.449	0.623	3.270	7.834

UniAud is superior than theirs across different $m = n$.

In general, for both tasks, ϵ_L increases and then decreases when m gets larger for all methods, which reflecting the expected trade-off controlled by m . The optimal choice of m for maximizing auditing performance is around $5 \times 10^2 - 2 \times 10^3$ across tasks in case I with $n = m$.

3) *Effectiveness of Orthogonal Features*: As shown in Figure 4 and Table IV, orthogonal features outperforms gaussian features across different audit data size, especially for a relatively small m or n . The critical point is attributed to the feature dimension that we set $d_x = 10^3$. Because the feature vectors of our orthogonal variant become approximated orthogonal when $m > d_x$, which explains the drop of ϵ_L in the left sub-figure for $m > 10^3$. Thus, it is expected that a larger d_x would result in a larger critical point of m for our orthogonal method.

4) *Effectiveness of Self-Comparison*: Above auditing improvement comes from our canary construction with uncorrelated data and the self-comparison framework. It should be noted that the self-comparison cannot stand without the uncorrelated canaries. So we now ablate effectiveness of self-comparison for further analysis.

As shown in Figure 5, for both $m = 2 \times 10^3$ and $m = 1 \times 10^4$, self-comparison significantly outperforms than the version without it, even with the same canary construction. By comparing the version w/o self-comparison with real-data best

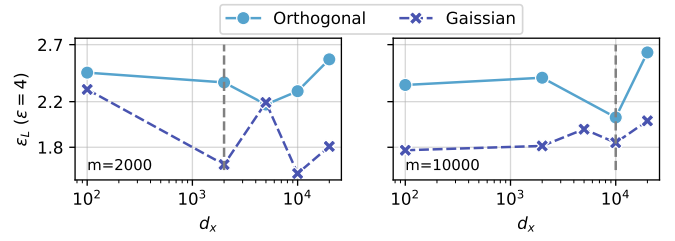


Fig. 6: Influence of feature dimension d_x .

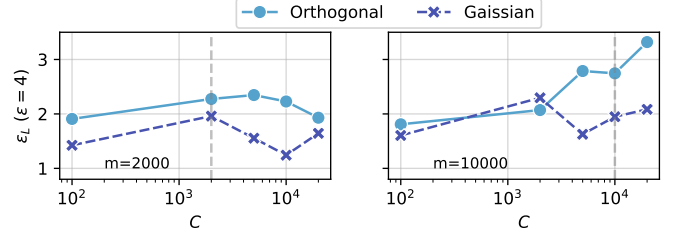


Fig. 7: Influence of the number of target encoding space C .

result (e.g., $\epsilon_L \approx 1.3$ for $\epsilon = 4$ [11], [29]), it shows a slight improvement due to the pure uncorrelated canary sample. This shows that the independent scoring f_{SCORE} and MIA decision f_{MIA} matters. While existing works on improving the canary only [29], [9], [11], UniAud improves more by consider to reduce the data dependency in **both** training and MIA inference.

5) *Influence of Canary Dimensions*: As observed in Figure 4, the input feature dimension matters especially for our orthogonal variant. Thus, we deeply analyze its influence in Figure 6. For two different $m \in \{2 \times 10^3, 1 \times 10^4\}$, we can observe that the estimated ϵ_L first decreases then increases with larger d_x after a critical point. The decrease is not monotonous and not as significant as the increase after the critical point. According to the gray dashed vertical line, it shows that the critical point comes after the $d_x \geq m$ (with $m = n$) for both our Gaussian and Orthogonal variants. Thus, it suggests that a safe choice of the dimension of input features $d_x \gg m$, aligning with a similar conclusion drawn in a recent auditing study [12].

6) *Influence of Encoding Space*: We find that the number of encoding space size C is vital hyper-parameters, because it depends how many m can be sampled uniquely, or in other words, how sparsity can the membership encoding space $\mathcal{X} \times \mathcal{Y}$ (or terms of $\mathcal{G} \times \mathcal{E}$ in Case II) is. In Figure 7, it is interesting to find that the optimal choice of C emerges after the critical point of $C = m$. Intuitively, it indicates that unique encoding is important for a good membership encoding performance.

While for relatively smaller m (e.g., 2×10^3) the ϵ_L slightly crease after the optimal C , ϵ_L continues to increase for larger $m = 1 \times 10^4$. The trade-off here is that the number of model parameters scales up with larger C , and the signal to noise ratio in DP training decreases with small m , which makes a larger model harder to converge. But it is not a issue when

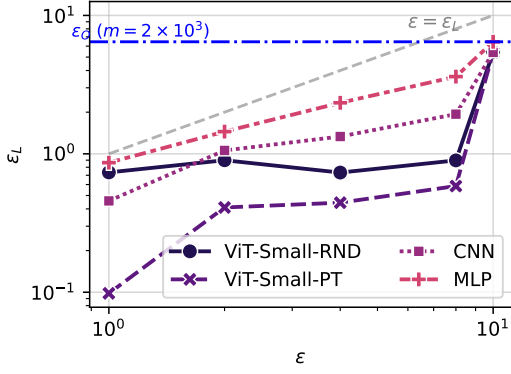


Fig. 8: Influence of the number of target encoding space C .

m is large enough (e.g., $m = 1 \times 10^4$), otherwise it can be solved by setting an even larger sub-sampling ratio in DP implementation.

7) *Influence of Pre-Training.*: It should be noted that our UniAud is independent of model architecture by design. While we assume auditor cannot craft the model initialization as *Poisoned*, the auditor can require the specific model architecture. For Case I where maxing auditing is the only criterion, it is meaningful to discuss what model should auditor designate for stronger auditing.

In Figure 8, we compare variants of UniAud on different model architectures, such as CNN with 8×10^5 , MLP with 3×10^7 and 2.2×10^7 parameters. Obviously, 2-layer MLP model with ReLU activations achieves the strongest auditing performance, probably because the large hidden states dimension $d_h = 10^5$. But the auditing strength does not simply scale with the trainable parameter size, we also notice that ViT-Small results in smaller ϵ_L . Comparing with training from scratch ViT-Small-RND with ViT-Small-PT, we find that random initialization benefits more for our synthetic canaries, especially for a small analytical privacy budget. This is because the mode initialization pre-trained on real-world dataset already finds a relatively flatten loss landscape, and it is more robust to noise. Otherwise, it requires more iterations for the pre-training model to fit the sample-specific noise in our canary. Thus, we suggest to use a random initialized neural network (such as 2-layer MLP with ReLU activations) in UniAud to maximize the auditing power.

C. Evaluation of Data-Dependent Auditing

In this section, we perform auditing for **Case II**, where both good model utility and auditing capability are criterion. For the baseline methods to achieve the best utility and auditing trade-off, we use the optimal hyper-parameters such as learning rate, batch size that we grid-searched for in Figure 1a.

1) *Trade-off for Image Classification*: As shown in Figure 9, we evaluate with two representative models: training from scratch for CNN and fine-tuning for pre-trained ViT-Small. We use a special instance of our random canary by repeating a randomly sampled pixel from a $3 \times 256 \times 256$

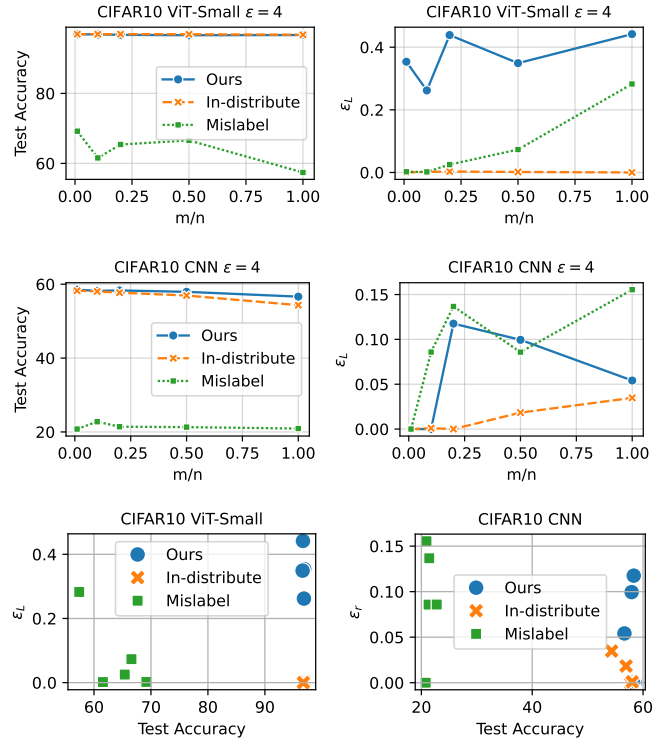


Fig. 9: Utility and Auditing Trade-off in Case II for Image Classification

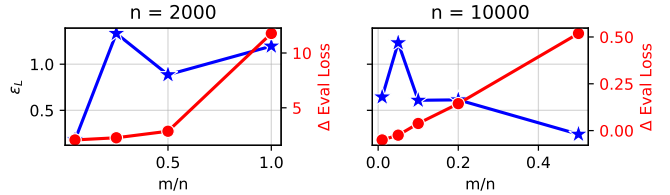


Fig. 10: Utility and Auditing Trade-off in Case II for language modeling on GPT-2. The Δ evaluation loss indicates the loss difference of UniAud++ with canaries minus the *In-distribution* baseline with all normal samples.

image, and apply it as a patch of dimension $d_x = 3 \times 10 \times 10$ to a real-world image. To reduce the required weights for the auditing head, we set $|\mathcal{E}| = 100$ in Algorithm 3.

It is clear that for both models, our method achieves the best utility-auditing trade-off than other $O(1)$ variants. For pre-trained models such as ViT, we notice that UniAud++ achieves the best of the word: near zero utility loss compared to $O(1)$ with original samples as canary; while large improvement on ϵ_L . While *Mislabeled* has better ϵ_L compared to *In-Distribution*, it significantly hurts the utility even when the ratio of mislabel m/n is small. The success of UniAud++ is attributed to our insight in Section IV-C of decoupling training objective, which disentangles the implicit conflict between utility and auditing.

It is also interesting to notice that for CNN models trained

from random initialization, UniAud++ can even outperform *In-Distribution* baseline that use normal samples as canary, while our ϵ_L is comparable to another baseline *Mislabeled* which absolutely abandons the utility. The key reason of why we have better utility compared *In-Distribution* is attributed to our Algorithm 2 which leverages m samples in training. By contrast, even though the canaries in *In-Distribution* are clean samples, the Bernoulli sampling in the original O(1) framework results in a loss of around $m/2$ normal samples for improving utility. And the loss on the amount of training data scales up with a larger m/n , as shown in Figure 9.

We do not consider comparison with *Poisoned* in case II because it requires the auditor to have a public in-distribution auxiliary dataset and craft the model initialization by warming up on the auxiliary dataset. Additionally, it is unfair to compare a pure private-preserving training with public-assisted training.

2) *Trade-off for Language Modeling*: We also analyze the utility and privacy trade-off in language modeling for case II. It should be noted that there is a distinguished difference between image classification task and language tasks in case II. In image classification task, the utility and auditing power conflicts because audit sample’s label share the same space as the normal non-audit samples. As an example, *Mislabeled* crafts the true label while still shares the space \mathcal{Y} . While for language models, when the training sample only includes augmented new tokens, although parameters of intermediate layers might change, the normal token embeddings do not change. Thus, as shown in Figure 10, when the total amount of training data n is relatively large, the utility degradation compared to *In-Distribution* with all normal samples is trivial.

In general, we observe a consistent trend for both $n = \{2 \times 10^3, 1 \times 10^4\}$ that the auditing capability (i.e., ϵ_L) first increases then decreases. The increase is caused with a higher statistical power limit on ϵ_O ; while the decrease is attributed to the degradation on the membership encoding sparsity which roughly scale according to $m/\sqrt{|V_{\text{new}}|}$.

IX. DISCUSSION AND BROADER IMPACTS

Towards practical auditing tools, our work builds on the efficient O(1) auditing framework under a realistic black-box access assumption. Although it is well known that combining these two settings makes auditing more challenging [11], [9], we believe it is necessary to position DP auditing as a seamlessly integrated component of machine learning pipeline – strengthening data governance and enabling continuous privacy monitoring, especially in sensitive domains such as healthcare analytics and government AI systems where theoretical guarantees demand empirical validation.

Unlike prior efforts that tighten theoretical auditing bounds via alternative confidence intervals [18], [41] or novel DP definitions [13], [12], our study offers fresh insights into black-box O(1) DP auditing by systematically reframing *the goal of DP auditing* from the perspective of *efficiency-utility-auditing* trade-off space—a critical dimension that has been largely overlooked [8], [10], [11].

By revealing the fundamental challenge of O(1) auditing, we categorize the auditing goals as: 1) Data-independent auditing, which targets the DP implementation itself; 2) Data-dependent auditing, which evaluates the empirical privacy risk of a specific training run on a given dataset and model. Our problem formulation opens two pivotal research directions:

- 1) Development of better membership encoding tasks that can maximize auditing tightness for the data-independent DP implementation. It is critical for monitoring the commercial Machine-Learning-as-a-Service compliance to privacy guidelines. More specifically, based on our formulation of membership encoding loss, future works are encouraged to design better canary data construction and model architecture for facilitating a tighter auditing.
- 2) Co-design of non-intrusive audit sub-tasks that intrinsically enhance membership signals *during* primary model training without hurting main task utility. This auditing-along-training framework is useful especially for the scenario when data-dependent risk is the target or when a single DP training is expensive. Based on our key insight of decoupling training objective, better auditing sub-task helps to further mitigate the fundamental conflict of utility and auditing.

Our methodology bridges the critical gap between formal privacy analysis and practical deployment constraints, establishing new pipelines for trustworthy data management in the widely available black-box learning scenarios. While we evaluate on image classification and language modeling, we anticipate extending both auditing frameworks to other tasks—such as image generation—via tailored membership encoding strategies and sub-task designs.

X. CONCLUSIONS

In this work, we rethink the fundamental challenge of black-box O(1) auditing and formulate two problems: data-independent and data-dependent auditing. We identify two bottlenecks: inevitable data dependency and implicit conflict between utility and auditing, which either remain as open problems or are ignored by existing works. To mitigate data dependency, we propose UniAud with a novel canary construction using uncorrelated features and labels, formulating auditing as a membership encoding task with a pair-matching objective. To further reduce dependency among audit samples during the membership inference stage, we introduce a self-comparison framework. We extend UniAud to UniAud++ for data-dependent scenarios, where audit results are estimated alongside models with acceptable utility, framing it as a novel multi-task problem by decoupling training objectives. Evaluation across image classification and language modeling tasks shows that our black-box O(1) data-independent auditing outperforms state-of-the-art approaches while saving thousands of training runs, and the data-dependent auditing achieves a better auditing-utility trade-off.

REFERENCES

- [1] L. Zhu, Z. Liu, and S. Han, “Deep leakage from gradients,” *Advances in neural information processing systems*, vol. 32, 2019.
- [2] N. Carlini, S. Chien, M. Nasr, S. Song, A. Terzis, and F. Tramer, “Membership inference attacks from first principles,” in *2022 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2022, pp. 1897–1914.
- [3] N. Carlini, F. Tramer, E. Wallace, M. Jagielski, A. Herbert-Voss, K. Lee, A. Roberts, T. Brown, D. Song, U. Erlingsson *et al.*, “Extracting training data from large language models,” in *30th USENIX Security Symposium (USENIX Security 21)*, 2021, pp. 2633–2650.
- [4] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang, “Deep learning with differential privacy,” in *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, 2016, pp. 308–318.
- [5] C. Dwork, F. McSherry, K. Nissim, and A. Smith, “Calibrating noise to sensitivity in private data analysis,” *Journal of Privacy and Confidentiality*, vol. 7, no. 3, pp. 17–51, 2016.
- [6] R. Shokri, M. Stronati, C. Song, and V. Shmatikov, “Membership inference attacks against machine learning models,” in *2017 IEEE symposium on security and privacy (SP)*. IEEE, 2017, pp. 3–18.
- [7] F. Tramer, A. Terzis, T. Steinke, S. Song, M. Jagielski, and N. Carlini, “Debugging differential privacy: A case study for privacy auditing,” *arXiv preprint arXiv:2202.12219*, 2022.
- [8] M. Jagielski, J. Ullman, and A. Oprea, “Auditing differentially private machine learning: How private is private sgd?” *Advances in Neural Information Processing Systems*, vol. 33, pp. 22 205–22 216, 2020.
- [9] M. Nasr, S. Song, A. Thakurta, N. Papernot, and N. Carlini, “Adversary instantiation: Lower bounds for differentially private machine learning,” in *2021 IEEE Symposium on security and privacy (SP)*. IEEE, 2021, pp. 866–882.
- [10] M. Nasr, J. Hayes, T. Steinke, B. Balle, F. Tramèr, M. Jagielski, N. Carlini, and A. Terzis, “Tight auditing of differentially private machine learning,” in *32nd USENIX Security Symposium (USENIX Security 23)*, 2023, pp. 1631–1648.
- [11] T. Steinke, M. Nasr, and M. Jagielski, “Privacy auditing with one (1) training run,” *Advances in Neural Information Processing Systems*, vol. 36, pp. 49 268–49 280, 2023.
- [12] Z. Xiang, T. Wang, and D. Wang, “Privacy audit as bits transmission:(im) possibilities for audit by one run,” *arXiv preprint arXiv:2501.17750*, 2025.
- [13] S. Mahloujifar, L. Melis, and K. Chaudhuri, “Auditing f -differential privacy in one run,” *arXiv preprint arXiv:2410.22235*, 2024.
- [14] M. S. Muthu Selva Annamalai and E. De Cristofaro, “Nearly tight black-box auditing of differentially private machine learning,” *Advances in Neural Information Processing Systems*, vol. 37, pp. 131 482–131 502, 2024.
- [15] A. Panda, T. Wu, J. Wang, and P. Mittal, “Differentially private in-context learning,” in *The 61st Annual Meeting of the Association for Computational Linguistics*, 2023.
- [16] B. Jayaraman and D. Evans, “Evaluating differentially private machine learning in practice,” in *28th USENIX Security Symposium (USENIX Security 19)*, 2019, pp. 1895–1912.
- [17] S. Maddock, A. Sablayrolles, and P. Stock, “Canife: Crafting canaries for empirical privacy measurement in federated learning,” *arXiv preprint arXiv:2210.02912*, 2022.
- [18] F. Lu, J. Munoz, M. Fuchs, T. LeBlond, E. Zaresky-Williams, E. Raff, F. Ferraro, and B. Testa, “A general framework for auditing differentially private machine learning,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 4165–4176, 2022.
- [19] S. Zanella-Béguelin, L. Wutschitz, S. Tople, A. Salem, V. Rühle, A. Paverd, M. Naseri, B. Köpf, and D. Jones, “Bayesian estimation of differential privacy,” in *International Conference on Machine Learning*. PMLR, 2023, pp. 40 624–40 636.
- [20] G. Andrew, P. Kairouz, S. Oh, A. Oprea, H. B. McMahan, and V. M. Suriyakumar, “One-shot empirical privacy estimation for federated learning,” *arXiv preprint arXiv:2302.03098*, 2023.
- [21] K. Pillutla, G. Andrew, P. Kairouz, H. B. McMahan, A. Oprea, and S. Oh, “Unleashing the power of randomization in auditing differentially private ml,” *Advances in Neural Information Processing Systems*, vol. 36, pp. 66 201–66 238, 2023.
- [22] R. Liu, T. Wang, Y. Cao, and L. Xiong, “Precurious: How innocent pre-trained language models turn into privacy traps,” in *Proceedings of the 2024 ACM SIGSAC Conference on Computer and Communications Security*, 2024.
- [23] Y. Wen, L. Marchyok, S. Hong, J. Geiping, T. Goldstein, and N. Carlini, “Privacy backdoors: Enhancing membership inference through poisoning pre-trained models,” *arXiv preprint arXiv:2404.01231*, 2024.
- [24] Z. Huang, N. Z. Gong, and M. K. Reiter, “A general framework for data-use auditing of ml models,” in *Proceedings of the 2024 on ACM SIGSAC Conference on Computer and Communications Security*, 2024, pp. 1300–1314.
- [25] Y. Tong, J. Ye, S. Zarifzadeh, and R. Shokri, “How much of my dataset did you use? quantitative data usage inference in machine learning,” in *The Thirteenth International Conference on Learning Representations*.
- [26] C. Song and R. Shokri, “Membership encoding for deep learning,” in *Proceedings of the 15th ACM Asia Conference on Computer and Communications Security*, 2020, pp. 344–356.
- [27] Z. Chen and K. Pattabiraman, “A method to facilitate membership inference attacks in deep learning models,” *arXiv preprint arXiv:2407.01919*, 2024.
- [28] P. Kairouz, S. Oh, and P. Viswanath, “The composition theorem for differential privacy,” in *International conference on machine learning*. PMLR, 2015, pp. 1376–1385.
- [29] A. Panda, X. Tang, M. Nasr, C. A. Choquette-Choo, and P. Mittal, “Privacy auditing of large language models,” *arXiv preprint arXiv:2503.06808*, 2025.
- [30] S. Wu, H. Zhang, G. Valiant, and C. Ré, “On the generalization effects of linear transformations in data augmentation,” in *International conference on machine learning*. PMLR, 2020, pp. 10 410–10 420.
- [31] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals, “Understanding deep learning requires rethinking generalization,” *arXiv preprint arXiv:1611.03530*, 2016.
- [32] M. Bellare and P. Rogaway, “Introduction to modern cryptography,” *Lecture Notes*, 2001.
- [33] A. Panda, X. Tang, S. Mahloujifar, V. Schwag, and P. Mittal, “A new linear scaling rule for private adaptive hyperparameter optimization,” in *International Conference on Machine Learning*. PMLR, 2024, pp. 39 364–39 399.
- [34] N. Papernot, A. Thakurta, S. Song, S. Chien, and Ú. Erlingsson, “Tempered sigmoid activations for deep learning with differential privacy,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 10, 2021, pp. 9312–9321.
- [35] F. Dörmann, O. Frisk, L. N. Andersen, and C. F. Pedersen, “Not all noise is accounted equally: How differentially private learning benefits from large sampling rates,” in *2021 IEEE 31st International Workshop on Machine Learning for Signal Processing (MLSP)*. IEEE, 2021, pp. 1–6.
- [36] Z. Bu, Y.-X. Wang, S. Zha, and G. Karypis, “Automatic clipping: Differentially private deep learning made easier and stronger,” *Advances in Neural Information Processing Systems*, vol. 36, 2023.
- [37] S. De, L. Berrada, J. Hayes, S. L. Smith, and B. Balle, “Unlocking high-accuracy differentially private image classification through scale,” *arXiv preprint arXiv:2204.13650*, 2022.
- [38] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, “An image is worth 16x16 words: Transformers for image recognition at scale,” *arXiv preprint arXiv:2010.11929*, 2020.
- [39] X. Li, F. Tramer, P. Liang, and T. Hashimoto, “Large language models can be strong differentially private learners,” *arXiv preprint arXiv:2110.05679*, 2021.
- [40] D. Yu, S. Naik, A. Backurs, S. Gopi, H. A. Inan, G. Kamath, J. Kulkarni, Y. T. Lee, A. Manoel, L. Wutschitz *et al.*, “Differentially private fine-tuning of language models,” *arXiv preprint arXiv:2110.06500*, 2021.
- [41] G. Zeng, W. Yang, Z. Ju, Y. Yang, S. Wang, R. Zhang, M. Zhou, J. Zeng, X. Dong, R. Zhang *et al.*, “Meddialog: Large-scale medical dialogue datasets,” in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2020.