IDENTIFYING OPTIMAL OUTPUT SETS FOR DIFFEREN TIAL PRIVACY AUDITING

Anonymous authors

Paper under double-blind review

ABSTRACT

Differential privacy limits an algorithm's privacy loss, defined as the maximum influence any individual data record can have on the probability of observing any possible output. Privacy auditing identifies the worst-case input datasets and output event sets that empirically maximize privacy loss, providing statistical lower bounds to evaluate the tightness of an algorithm's differential privacy guarantees. However, current auditing methods often depend on heuristic or arbitrary selections of output event sets, leading to weak lower bounds. We address this critical gap by introducing a novel framework to compute the optimal output event set that maximizes the privacy loss lower bound in auditing. Our algorithm efficiently computes this optimal set when closed-form output distributions are available and approximates it using empirical samples when they are not. Through extensive experiments on both synthetic and real-world datasets, we demonstrate that our method consistently tightens privacy lower bounds for auditing differential privacy mechanisms and black-box DP-SGD training. Our approach outperforms existing auditing techniques, providing a more accurate analysis of differentiallyprivate algorithms.

025 026 027

028 029

035

004

010 011

012

013

014

015

016

017

018

019

021

023

1 INTRODUCTION

Differential privacy (DP) (Dwork et al., 2014) is regarded as the gold standard for quantitatively assessing an algorithm's privacy guarantee. It ensures that the presence or absence of any record will not be heavily revealed from the outcome of the mechanism. Formally, consider a randomized algorithm $M : \mathcal{D} \to \Theta$. If for any neighbouring datasets $D, D' \in \mathcal{D}$ that differ by only one entry, i.e., $D' = D \cup \{x\}$ or $D' = D \setminus \{x\}$, and for any measurable (output event) sets \mathcal{O} , it satisfies that

$$\mathbb{P}[M(D) \in \mathcal{O}] \le e^{\epsilon} \cdot \mathbb{P}[M(D') \in \mathcal{O}] + \delta,$$

then the algorithm M is said to satisfy (ϵ, δ) -DP. Differential privacy mathematically certifies that the information leakage of a randomized algorithm is bounded (by ε and δ) even under the *worst-case* datasets D, D' and the *worst-case* output event set O.

Computing the exact differential privacy loss ε and δ for an algorithm requires tracking its output 040 distribution over all possible input datasets and output event sets, which is challenging for com-041 plex algorithm (e.g., iterative learning algorithms with a large number of non-linear updates). To 042 circumvent this challenge, existing privacy analysis for machine learning algorithms (Abadi et al., 043 2016; Kairouz et al., 2015; Mironov, 2017) often make simplifying assumptions that the adversary 044 observes not only the final output, but also other information (such as intermediate updates) in a 045 learning algorithm, and prove upper bounds for ε and δ instead. Understanding whether the DP up-046 per bound is correct and tight is critical, since underestimating privacy upper bounds leads to privacy 047 leaks (Tramer et al., 2022), while any underestimation results in unnecessarily large modifications 048 to the learning algorithms which decreases the utility of the trained model. To measure the tightness of the DP upper bound, it is crucial to identify the worst-case datasets and output event sets where the privacy loss is maximized. The success of inference attack under such worst-case scenarios then 051 provides a statistical lower bound on the exact differential privacy loss of the algorithm. If the lower bound is close to the DP upper bound, then the privacy analysis is tight. This process is known as 052 privacy auditing (Jagielski et al., 2020; Nasr et al., 2021; Zanella-Beguelin et al., 2023; Nasr et al., 2023; Steinke et al., 2023).

084

085

090

091 092

094

095

096

097

098

099

054 Existing privacy auditing literature largely focuses on designing the worst-case input datasets D, D'via poisoning attacks (Jagielski et al., 2020; Nasr et al., 2021) or via searching for the data with 056 the highest risks against inference attacks in the training dataset via computationally expensive inference attacks (Ye et al., 2022; Zanella-Beguelin et al., 2023; Aerni et al., 2024). By contrast, the 058 question of how to choose the worst-case output event set \mathcal{O} to enable a better auditing lower bound is largely underexplored. This is due to several reasons. (1) Standard learning algorithms consist of a large number of iterative non-linear updates, which makes it **challenging to track the complex** 060 distribution of the final output in a *closed-form* manner. Consequently, it is *infeasible* to perform 061 *exact* optimization of output event set for complex iterative learning algorithms. (2) If we relax 062 ourselves to approximate optimization of output set, prior works (Bichsel et al., 2021; Lu et al., 063 2022) have proposed several approaches based on approximating the output distributions via their 064 empirical Monte Carlo samples. However, such approximations require a large number of Monte 065 Carlo samples to be accurate, which requires computationally expensive (re)runs of the training al-066 gorithm. Alternatively, given a small finite number of Monte Carlo samples, such methods suffer 067 from non-negligible approximation error, which in turn affects the auditing performance. Thus, it 068 is crucial to take the approximation error into account when formulating the objective of maximal 069 privacy loss in privacy auditing. However, to the best of our knowledge, there is no existing consensus regarding the optimal way to incorporate the approximation error into the outpupt set selection objective for privacy auditing. (3) Prior works (Steinke et al., 2023; Bichsel et al., 2021; 071 Lu et al., 2022) have empirically observed that the effect of output event set choice on the privacy 072 auditing power is distribution-dependent. For example, in certain mechanisms like randomized re-073 sponse, the choice of output event set does not significantly affect the audited lower bound ((Steinke 074 et al., 2023, Figure 10)). In contrast, for some other mechanisms (such as the Laplace mechanism), 075 the choice of output event set substantially affects the audited lower bound Bichsel et al. (2021); Lu 076 et al. (2022). Despite these observations, a comprehensive understanding of whether and when 077 optimizing the output event set yields a gain remains elusive (i.e., how does the amount of gain 078 change for different mechanisms). 079

In this paper, we aim to address the question of optimal output set selection for privacy auditing.
 Specifically, we ask the following research questions.

- 1. What is the optimal output event set for auditing differentially private mechanisms? How does the optimal output event set change for different mechanisms?
- 2. How much gain for privacy auditing can we achieve by optimizing the output set? How does the amount of gain depend on the mechanism and the number of auditing examples?
- 3. How to compute the optimal output event set when only given empirical output samples, rather than closed-form output distributions? Does this enable tighter auditing lower bounds for practical DP learning mechanisms, such as black-box DP-SGD?

Summary of Results. Our results for addressing the above research questions are as follows.

- We proposed a novel framework for choosing the optimal output set in privacy auditing, that incorporates the error of sampling auditing examples from output distributions. Specifically, our optimization objective (4) captures the expectation of advantage-based auditing lower bound (Steinke et al., 2023; Jayaraman & Evans, 2019) over random sampling.
- We prove that our proposed optimization problem incurs closed-form solutions, that consist of output sets constructed by thresholding the likelihood ratio between member and non-member score distributions in privacy auditing (Definition 4.2). Importantly, we prove (Theorem 4.3) that our constructed output set enables the *maximum auditing objective* among all possible choices of output event set, i.e., is optimal.
- We numerically validate that for auditing DP mechanisms with closed-form output distributions, using the optimal output set *consistently* enables tighter lower bounds than prior methods. Our improvement is especially significant under a small number of auditing examples (Figure 2), or when the output distribution is asymmetric or multi-modal (Figure 1).
- We also extend our framework to algorithms *without* closed-form output distributions (Section 4.3), and approximately compute the optimal output set by thresholding the KDE estimation of likelihood ratio from empirical output samples. We experimentally demonstrate

(Section 6) that our algorithm is effective and enables tighter privacy auditing of black-box DP SGD training on synthetic (canary) and real-world (CIFAR-10) datasets.

111

108

109

110

112

2 RELATED WORKS

113 Score-based membership inference attacks and connections to output set selection. Membership 114 inference attacks (MIAs) aim to infer whether a given data record was used for training a given target 115 model. Intuitively, a higher performance of MIA indicates higher privacy risks, and thus MIAs serve 116 as crucial tools for auditing DP mechanisms (Shokri et al., 2017). MIAs in the literature generally 117 follow a fixed template: compute a membership inference attack score (such as loss (Sablayrolles 118 et al., 2019; Yeom et al., 2018; Ye et al., 2022), confidence (Salem et al., 2018), entropy (Song & 119 Mittal, 2021), loss trajectory (Liu et al., 2022) and likelihood ratio (Carlini et al., 2022)) on the 120 target model and target data, and then perform attack via thresholding the MIA score. That is, 121 predict member if the MIA score is above the threshold, and predict non-member if the MIA score is below the threshold. To this end, performing score-based membership inference attacks (on the 122 output distributions in privacy auditing) is *equivalent* to heuristically selecting the output event set 123 to be the interval (Z, ∞) by a threshold Z, as utilized in the auditing experiment in Jagielski et al. 124 (2020). This strategy, as we will show in Section 4, 5 and 6, is strictly suboptimal for certain DP 125 mechanisms compared to our proposed optimization-based output event set selection. 126

127 Optimization of output event set in privacy auditing. Bichsel et al. (2021) optimizes the posterior over output event set, while ensuring that the probability of output event set is larger than a heuristi-128 cally selected threshold (e.g., 0.05). Lu et al. (2022) further trains a ML model to learn the posterior 129 of training data given trained model, and selects the output set that maximizes the posterior. Steinke 130 et al. (2023) considers a series of output event sets constructed by a fixed strategy. Their constructed 131 output event set takes the form of two intervals $(-\infty, c_{-})$ and (c_{+}, ∞) , where c_{-} and c_{+} are real-132 valued thresholds for the MIA scores. They then search over the entire output domain (i.e., the set 133 of values for discrete Monte Carlo samples) to identify a choice of c_{-} and c_{+} that enables maxi-134 mum empirical auditing performance. However, prior works do not reach a consensus regarding the 135 optimal output set for privacy auditing. By contrast, our work proposes a novel objective (4) that 136 captures the expectation of auditing lower bound and serves as the key ingredient for us to prove an 137 optimality guarantee for selecting output set in privacy auditing given finite auditing samples.

Evaluating worst-case success of membership inference attacks. A connected direction of research is on measuring privacy risk of machine learning algorithms by the success of membership inference attack over a small set of (worst-case) data samples, rather than the entire dataset. To this end, Carlini et al. (2022) propose to evaluate the true positive rate (TPR) at a low false positive rate (FPR) as an indicator of privacy loss as worst-case data records. This is equivalent to evaluating MIA only on a selected output set at the bottom tail of the MIA score distributions on all training and test data (as this constrains the FPR to be small), rather than the whole output domain.

145 146 147

3 PROBLEM STATEMENT

148 149 149 149 149 150 151 152 153 Let $\mathcal{T}: D \mapsto \theta$ be an (ε, δ) -differentially private algorithm, that maps an input dataset D to a trained model with parameters θ . Let $Score(\theta, x)$ be a MIA score function, that maps the algorithm's output model parameters θ and any data record x to a real-valued score o. An auditing experiment \mathcal{E} takes in an randomized algorithm \mathcal{T} , a confidence tolerance β , and an approximate differential privacy tolerance δ as inputs, and returns a high-confidence statistical lower bound $\hat{\varepsilon} := \mathcal{E}(\mathcal{T})$ for the ground-truth differential privacy parameter ε that satisfies

154

$$\Pr[\hat{\varepsilon} \coloneqq \mathcal{E}(\mathcal{T}) < \varepsilon] \ge 1 - \beta \tag{1}$$

where the probability is over randomness of the auditing experiment \mathcal{E} and the training algorithm \mathcal{T} . Generally speaking, existing auditing experiments (Jagielski et al., 2020; Nasr et al., 2021; 2023; Zanella-Beguelin et al., 2023; Bichsel et al., 2021; Lu et al., 2022) (see Appendix A and Table 1 for more details) can be decomposed into the following three components.

- 159 160 161
- 1. Compute a set of member scores S_+ and non-member scores S_- . This involves sampling the input dataset, training output models, and computing MIA scores on the trained model(s) and its member and non-member data record(s) respectively.

- 2. Subselect the scores by a output set $\mathcal{O} \subseteq \mathbb{R}$. Intuitively, this is to increase the distinguishability between subselected member and non-member scores (that lie in \mathcal{O}), so as to enable a higher auditing lower bound.
 - 3. Design auditing function L such that $L(S_+, S_-, \mathcal{O}; \delta, \beta) \leq \varepsilon$ with high probability. Typically, L is a function of the performance of membership inference in auditing. ¹ To prove L satisfies (1), prior works generally rely on the constraints of (ε, δ) -DP on the power of binary membership hypothesis test (Wasserman & Zhou, 2010; Kairouz et al., 2015).

Our objective We aim to design a general framework for computing the optimal output set, so as to enable the highest possible lower bound $\hat{\varepsilon}$ in privacy auditing. Formally speaking, for any fixed distribution of member scores S_+ and non-member scores S_- (specified by the data sampling and training protocols in the auditing experiment), any approximate DP tolerance δ , any confidence tolerance β , and any auditing function L, we aim to solve the below optimization problem.

162

163

164

166

167

 $\arg\max_{\mathcal{O}} \hat{\varepsilon} = L(S_+, S_-, \mathcal{O}; \delta, \beta) \tag{2}$

176 By fixing S_+, S_-, δ, β , and L, we are essentially fixing the first and third components of the auditing 177 experiment. For simplicity of presentation, and for fair comparison with prior works (Bichsel et al., 178 2021; Lu et al., 2022) that solely focus on pure DP settings, in the paper, we only optimize and 179 report auditing lower bounds for $\delta = 0$ and $\beta = 0.05$. To compute the member and non-member scores S_+, S_- , we focus on two representative ways in the literature (Jagielski et al., 2020; Steinke 181 et al., 2023) as summarized by Table 1. For auditing function, we will use generalized variants (Proposition 4.1) of the advantage-based auditing function in (Steinke et al., 2023, Theorem 5.2).² 182 However, we emphasize that we aim to design a general framework for optimizing the output set that 183 can be combined with any member and non-member scores and auditing functions in the literature. 184

185 186

187

4 COMPUTING THE OPTIMAL OUTPUT SET FOR PRIVACY AUDITING

In this section, we introduce our optimization objective for selecting the output set to enable higher advantage-based privacy auditing lower bound. We focus on the setting where the member and nonmember scores S_+ and S_- are i.i.d. samples from two output distributions p and q respectively. This is indeed the case for prototypical auditing experiments (Jagielski et al., 2020; Nasr et al., 2021) – see Appendix A and Table 1 for details. We then analyze the structure of the optimal output set and propose an algorithm to approximate the optimal output set from empirical samples.

194 195

4.1 OPTIMIZATION OBJECTIVE FOR OUTPUT SET SELECTION

Given black-box access to i.i.d. samples from the output distributions of the DP mechanism, we first prove a new auditing function that has explicit dependence on the choice of output set \mathcal{O} as follows.

Proposition 4.1 (Our Output-set-dependent Auditing Function). Let p and q be two probability distributions over \mathbb{R} that satisfies $e^{-\varepsilon} \leq \frac{p(o)}{q(o)} \leq e^{\varepsilon}$ for any $o \in \mathbb{R}$. Let $\beta = 0.05$ be the confidence tolerance. Let $\mathcal{O} \subset \mathbb{R}$ be a fixed output set. Let S_+ be m i.i.d. samples from p, and let S_- be mi.i.d. samples from q. Then $\hat{\varepsilon}(S_+, S_-, \mathcal{O}; p, q)$ defined as follows is a valid auditing function.

203 204

$$\hat{\varepsilon}(S_+, S_-, \mathcal{O}; p, q) = \phi\left(\frac{\int_{o \in \mathcal{O}} \max\{p(o), q(o)\}do}{p(\mathcal{O}) + q(\mathcal{O})} - \sqrt{\frac{2 \cdot \ln(2/\beta)}{|S_+ \cap \mathcal{O}| + |S_- \cap \mathcal{O}|}}\right)$$
(3)

209

210

211

where $\phi(y) = \log\left(\frac{y}{1-y}\right)$ is the logit function. That is, $\Pr[\hat{\varepsilon}(S_+, S_-, \mathcal{O}; p, q) \le \varepsilon] \ge 1 - \beta$.

Proof. The proof is in Appendix B.3. This generalizes prior auditing functions Yeom et al. (2018); Steinke et al. (2023) from specific output sets to any \mathcal{O} , and enjoys the benefit of explicit dependence on the output set – see Remark B.4 and B.5 for details.

²¹² ¹Other auditing functions include the ones based on attribute inference Guo et al. (2023), reconstruction Guo ²¹³ et al. (2022); Balle et al. (2022), and empirical divergence estimation Domingo-Enrich & Mroueh (2022); Kong ²¹⁴ et al. (2024). It is interesting to invetigate the effect of output set optimization for them as future works. ²¹⁵ ²That is, the design of *L* is based on bounding the advantage of membership inference with the DP guarantee

²That is, the design of L is based on bounding the advantage of membership inference with the DP guarantee (ε, δ) . Extending our framework to other performance metrics of MIA is an interesting future work.

By taking expectation of the denominator $|S_+ \cap \mathcal{O}| + |S_- \cap \mathcal{O}|$ in (3) across random sampling of S_+ and S_- , we obtain an objective for selecting the optimal output set \mathcal{O} in privacy auditing.

$$\arg\max_{\mathcal{O}} \hat{\varepsilon}(\mathcal{O}; p, q) \coloneqq \phi \left(\underbrace{\frac{\int_{o \in \mathcal{O}} \max\{p(o), q(o)\} do}{p(\mathcal{O}) + q(\mathcal{O})}}_{\text{inference accuracy}} - \underbrace{\sqrt{\frac{2 \cdot \ln(1/\beta)}{m \cdot (p(\mathcal{O}) + q(\mathcal{O}))}}}_{\text{sampling error}} \right)$$
(4)

where $\phi(y) = \log\left(\frac{y}{1-y}\right)$. Note that the objective (4) is specified by distributions p and q, indicating that the optimal \mathcal{O} is mechanism-dependent. We now analyze the structure of the optimal output set.

4.2 IDENTIFYING OPTIMAL OUTPUT SET FOR AUDITING

We first define a series of output set \mathcal{O}_{τ} by thresholding the likelihood ratio between distributions. **Definition 4.2** (τ -log-likelihood-ratio-set). Let p and q be two continuous distributions over \mathbb{R} , indicating member and non-member score distributions in the auditing experiment respectively. For each $\tau > 0$, we define the τ -log-likelihood-ratio-set for p and q as the following set \mathcal{O}_{τ} .

$$\mathcal{O}_{\tau} = \left\{ x \in \mathbb{R} : \left| \ln \left(\frac{p(x)}{q(x)} \right) \right| \ge \tau \right\}$$
(5)

We now prove the optimality guarantee for \mathcal{O}_{τ} , i.e., \mathcal{O}_{τ} enables the largest lower bound among all possible output sets S that satisfy the same size constraint $p(\mathcal{O}) + q(\mathcal{O}) = p(\mathcal{O}_{\tau}) + q(\mathcal{O}_{\tau})$.

Theorem 4.3. Let p and q be two continuous distributions over \mathbb{R} . Let $\hat{\varepsilon}(\cdot; p, q)$ be our auditing objective (4). Given any feasible output set \mathcal{O} , there exists $\tau \in \mathbb{R}$ such that $p(\mathcal{O}_{\tau}) + q(\mathcal{O}_{\tau}) = p(\mathcal{O}) + q(\mathcal{O})$, where \mathcal{O}_{τ} be the τ -log-likelihood-ratio-set (5) for p and q. Further, it satisfies that

$$\hat{\varepsilon}(\mathcal{O}_{\tau}; p, q) = \max_{\mathcal{O} \subseteq \mathbb{R}: p(\mathcal{O}) + q(\mathcal{O}) = p(\mathcal{O}_{\tau}) + q(\mathcal{O}_{\tau})} \hat{\varepsilon}(\mathcal{O}; p, q)$$
(6)

The proof technique (Appendix B.2) is similar to the Neyman-Pearson Lemma (Neyman & Pearson, 1933) (Remark B.6). Theorem 4.3 proves that for any output set \mathcal{O} , there exists a τ -log-likelihoodratio-set \mathcal{O}_{τ} that satisfies $p(\mathcal{O}_{\tau}) + q(\mathcal{O}_{\tau}) = p(\mathcal{O}) + q(\mathcal{O})$ such that \mathcal{O}_{τ} enables higher auditing lower bound objective (4) than \mathcal{O} . Thus the family of $\{\mathcal{O}_{\tau}\}_{\tau>0}$ contains the optimal output set.

4.3 APPROXIMATING THE OPTIMAL OUTPUT SET FROM EMPIRICAL SAMPLES

In practice, we are only given empirical samples from the output distributions p and q, rather than the closed-form densities. We now propose a method to approximate the optimal output set \mathcal{O}_{τ} from empirical samples. The idea is to estimate the probability densities of distributions p and qfrom their empirical samples (e.g., via kernel density estimation), and then compute the likelihood ratio function and its level set as optimal output set. See Algorithm 1 for the pseudocode.

Algorithm 1 Approximating the Optimal Optimal Output Set from Empirical Samples

255 **Require:** m_+ samples S_+ and m_- samples S_- from unknown distributions p, q respectively, with $m_+ + m_- = 2m$. DP tolerance δ . Confidence β . Auditing function $L(S_+, S_-, \mathcal{O}; \delta, \beta)$. 256 1: Kernel density estimation (KDE) from samples: $\hat{p} \leftarrow KDE(S_+), \hat{q} \leftarrow KDE(S_-)$ 257 2: Sort the empirical samples: $\tilde{o}_0 = -\infty$; $\tilde{o}_1 \le \cdots \le \tilde{o}_{2m} \leftarrow \operatorname{sort}(S_+ \cup S_-)$; $\tilde{o}_{2m+1} = +\infty$ 3: Estimate absolute log likelihood ratio: $\{\tau_0, \cdots, \tau_{2m}\} \coloneqq \left\{ \left| \log \frac{\int_{\tilde{o}_i}^{\tilde{o}_{i+1}} \hat{p}(o)do}{\int_{\tilde{o}_i}^{\tilde{o}_{i+1}} q(o)do} \right| \right\}_{i=0}^{2m}$ 258 259 260 261 4: for each level $\tau \in \{\tau_0, \cdots, \tau_{2m}\}$ do approximate the τ -log-likelihood-ratio set: $\mathcal{O}_{\tau} = \bigcup_{i=0:\tau_i > \tau}^{2m} (\tilde{o}_i, \tilde{o}_{i+1}).$ 262 5: 263 6: Search the optimal level $\hat{\tau} \coloneqq \arg \max_{\tau \in \{\tau_0, \dots, \tau_{2m}\}} L(S_+, S_-, \mathcal{O}_{\tau}; \delta, \beta).$ 264 7: **return** output set $\mathcal{O}_{\hat{\tau}}$.

265 266

219220221222223

224 225

226 227

228

229

230

231

236

237

238

239

240

241

246

253

254

267 On reducing the instability of KDE and likelihood ratio estimation Motivated by Bichsel et al. 268 (2021), we only perform likelihood ratio estimation on the regions where the estimated densities 269 are non-negligible (with probability density function $\hat{p}(o)$ and $\hat{q}(o)$ that is larger than 0.01). This is 269 because the likelihood ratio estimation involves division over probability density, and thus the error 269 can be arbitrarily large when the estimated densities are small. 270 **On searching for the optimal leve** $\hat{\tau}$ Theorem 4.3 proves that the family of τ -log-likelihood-ratio-271 set contains the optimal output set. To choose one single output set over the family of \mathcal{O}_{τ} , in Line 272 6 of Algorithm 1, we evaluate the auditing performance of \mathcal{O}_{τ} for different values τ and return *one* 273 output set with the best performance. This is a one-dimensional optimization problem over $\tau \in \mathbb{R}$, 274 which is significantly easier and incurs significantly less computation cost than the original output set optimization problem over all possible output sets $\mathcal{O} \subseteq \mathbb{R}$. Although this search for optimal 275 threshold requires additional output samples, it is a common practice in the literature (Bichsel et al., 276 2021; Lu et al., 2022; Zanella-Beguelin et al., 2023) to boost the auditing performance. 277

- 278
- 279 280

295 296 297

306

320 321

5 NUMERICAL EXPERIMENTS: ONE-EPOCH OF SHUFFLED NOISY SGD

281 In this section, we perform numerical experiments on two fundamental differentially private mech-282 anisms: Lapalace mechanism (which is pure DP), and Gaussian mechanism (which is approx-283 imate DP), that are building blocks for the iterative update in the celebrated DP-SGD learning algorithm (Abadi et al., 2016). We consider composition of both mechanisms for one shuffled epoch 284 of noisy stochastic gradient descent udpates (Bassily et al., 2014), under simple loss functions con-285 struction. We then investigate whether our output set optimization method enables tighter privacy 286 auditing lower bound than prior works (Jagielski et al., 2020; Bichsel et al., 2021; Lu et al., 2022; 287 Steinke et al., 2023), and how various factors affect the amount of gain. 288

Experiment setting Consider a simplified setting of learning on dataset $D = (x_1, x_2)$ with two records $x_1, x_2 \in \mathbb{R}$, where the loss function is $\ell_1(x, \theta) = \langle \theta, x \rangle$ for the first record, and is $\ell_2(x, \theta) = \frac{1}{2}(\theta - x)^2$ for the second record. Assuming different loss function for different records of the dataset is realistic, for example, in a multi-task learning setting, where the first record is used for learning task one (e.g., classify between cats and dogs) and the second record is used for learning task two (e.g., classify between sketch and photo). In terms of learning algorithm, we consider running one epoch of shuffled noisy stochastic gradient descent algorithm with the following updates:

$$\theta_{1} = \theta_{0} - \eta \left(\nabla \ell_{1}(\theta_{0}, x_{s_{1}}) + Z_{1} \right) \quad \text{where } Z_{1} \sim \text{Lap}(0, b_{1}) \text{ or } \mathcal{N}(0, \sigma_{1}^{2}) \\ \theta_{2} = \theta_{1} - \eta \left(\nabla \ell_{2}(\theta_{1}, x_{s_{2}}) + Z_{2} \right) \quad \text{where } Z_{2} \sim \text{Lap}(0, b_{0}.5) \text{ or } \mathcal{N}(0, \sigma_{2}^{2})$$
(7)

where θ_0 is the initialization parameters of the model (assumed to be zero for simplicity), η is the learning rate, b_1 , $b_0.5$, σ_1 and σ_2 are the noise scales, and $s \xleftarrow{\text{uniform}} \{(1,2),(2,1)\}$ is a randomly shuffled dataset order. Under the assumption of bounded data domain, the gradient of both loss functions $\ell_1(\theta, x)$ and $\ell_2(\theta, x)$ have finite sensitivity, thus the algorithm satisfies differential privacy by standard DP guarantees of Laplace mechanism (Dwork et al., 2006, Theorem 1) and Gaussian mechanism (Balle & Wang, 2018, Theorem 8). (See Appendix C for the details.) Conditioned on a fixed order *s*, the final output θ_2 of this mechanism follows a closed-form distribution as follows.

$$\theta|s \stackrel{d}{=} \theta_0 - \eta \cdot (1 - \eta) \cdot x_{s_1} + \eta \cdot x_{s_2} + Z \quad \text{where } Z = -\eta (1 - \eta) Z_1 - \eta Z_2 \tag{8}$$

For simplicity, assume that $\theta_0 = 0$ and $\eta = \frac{1}{2}$ (note that both terms do not affect the DP guarantee of the mechanism in Proposition C.1). By averaging over the two possible orders *s*, we could obtain the closed-form output distributions $p(\theta)$ and $q(\theta)$ for model θ trained on dataset $D = (x_1, x_2)$ and $D' = (x'_1, x'_2)$ respectively. For example, for Gaussian mechanism, when $Z_1 \sim \mathcal{N}(0, \sigma^2)$ and $Z_2 \sim \mathcal{N}(0, \sigma^2)$, we have that $p(\theta)$ and $q(\theta)$ follows the below mixture distributions.

$$p(\theta) = \frac{1}{2} \cdot \mathcal{N}\left(-\frac{1}{4}x_1 + \frac{1}{2}x_2, \frac{5}{16}\sigma^2\right) + \frac{1}{2} \cdot \mathcal{N}\left(-\frac{1}{4}x_2 + \frac{1}{2}x_1, \frac{5}{16}\sigma^2\right)$$
(9)

$$q(\theta) = \frac{1}{2} \cdot \mathcal{N}\left(-\frac{1}{4}x_1' + \frac{1}{2}x_2', \frac{5}{16}\sigma^2\right) + \frac{1}{2} \cdot \mathcal{N}\left(-\frac{1}{4}x_2' + \frac{1}{2}x_1', \frac{5}{16}\sigma^2\right)$$
(10)

Similarly, for Laplace mechanism, when $Z_1 = 0 \sim \text{Lap}(0,0)$ and $Z_2 \sim \text{Lap}(0,b)$, we have that $p(\theta)$ and $q(\theta)$ follows the below mixture distributions.

$$p(\theta) = \frac{1}{2} \cdot \operatorname{Lap}\left(-\frac{1}{4}x_1 + \frac{1}{2}x_2, \frac{1}{2}b\right) + \frac{1}{2} \cdot \operatorname{Lap}\left(-\frac{1}{4}x_2 + \frac{1}{2}x_1, \frac{1}{2}b\right)$$
(11)

322
323
$$q(\theta) = \frac{1}{2} \cdot \operatorname{Lap}\left(-\frac{1}{4}x_1' + \frac{1}{2}x_2', \frac{1}{2}b\right) + \frac{1}{2} \cdot \operatorname{Lap}\left(-\frac{1}{4}x_2' + \frac{1}{2}x_1', \frac{1}{2}b\right)$$
(12)

Equation (9), 10, 11 and 12 are the forms of output distribution p and q that we will use for auditing one-epoch shuffled noisy SGD in this section. That is, to compute member and non-member scores S_+ and S_- , we perform Monte Carlo sampling from the output distributions of the mechanism. This is equivalent to the auditing experiment via retraining in Jagielski et al. (2020).

Comparison Baselines For our method, we approximate the optimal output set by applying Algorithm 1 on top of black-box Monte Carlo samples from the output distributions((9), 10, 11 and 12) and auditing function add ref for i.i.d. samples given by prior works. To validate the optimality guarantee (Theorem 4.3), we compare with the following methods.

- 1. Whole domain: A naive strategy is to choose \mathcal{O} as \mathbb{R} , i.e., the whole score domain.
- 2. Top (bottom) score heuristic: (Steinke et al., 2023, Algorithm 1) uses a strategy of 'abstaining' by only performing MIAs on the top k^+ and bottom k_- scores in the auditing experiment, which is equivalent to the output set spanned by top α_+ -percentile and bottom α_- -percentile of the mixture of member and non-member scores. This selection intuitively restricts the guesses to scores that are most likely to be members or non-members, thus boosting the attack accuracy and the auditing lower bound. Experimentally, (Steinke et al., 2023, Figure 11) observes that the choice of k^+ and k_- has a significant impact on the auditing lower bound, indicating that this heuristic outperforms the naive baseline of using the whole score domain. In our experiments, we follow (Steinke et al., 2023) and bruteforce search over all possible values to find the best k_+ and k_- , and then report the auditing performance of the best k_+ and k_- in fresh evaluation trials.
- 3. **DP-sniper (Bichsel et al., 2021) and Lu et al. (2022):** Bichsel et al. (2021) propose to compute output set by thresholding the *posterior* probability of training dataset conditional on mechanism output. Lu et al. (2022) further makes the method more flexible by searching for the optimal threshold that enables the largest auditing lower bound. In experiments, we use the implementation ³ released by Bichsel et al. (2021) and additionally follow Lu et al. (2022) to search for the optimal threshold for fair comparison.

349 350

357

358

333

334

335

336

337

338

339

340

341

342 343

345

347

348

For all methods, we directly use the one-dimensional output of the DP mechanisms as the MIA score and subsequently perform the optimal advantage-based MIA with rejection region $A = \{o : \hat{p}(o) \ge \hat{q}(o)\}$ (benefited from the KDE estimators of the output distributions p and q). To evaluate the final auditing lower bound, we use the auditing function in (Steinke et al., 2023, Section D) for $\delta = 0$. For all methods, we use a separate set of output samples in evaluation, compared to the set of output samples used for searching for *one* optimal output set.

5.1 SHAPE OF OPTIMAL OUTPUT SET

We now investigate the shape of \mathcal{O}_{τ} for the shuffled noisy SGD under Laplace mechanism and Gaussian mechanism ((11), 12, 9, 10). Our observations are as follows.

Optimal set contains more than two intervals for multimodal distributions p and q In Figure 1, we observe that for mixture of Laplace or Gaussian distributions, the optimal output set contains more than two intervals. This means that the prior heuristic (Steinke et al., 2023) for selecting the top and bottom scores is strictly suboptimal. Specifically, the optimal output set \mathcal{O}_{τ} in Figure b contains the union of three disjoint intervals, one for the top-valued scores, one for bottomvalued scores, and the other for medium-valued scores. The deeper reason is that the likelihood ratio p(o)/q(o) is **not monotonic** in score domain o for distinguishing between general choices of distributions p and q that are multi-modal. See Appendix D for more examples of other distributions.

369 **Optimal set is asymmetric for general mechanisms** Motivated by Wilk's theorem (Wilks, 1938) 370 (which proves that the log-likelihood ratio statistics converges asymptotically to chi-squared dis-371 tribution under the null hypothesis in binary hypothesis testing), we additionally experiment on 372 Chi-squared distributions p and q in Appendix D and observe that the optimal output set is asym-373 metric. This is because the likelihood ratio p(o)/q(o) is **not symmetric** in score domain o due to 374 the asymmetric tails of chi-squared distributions p and q. This confirms the reason for the effectiveness of the "abstaining" strategy in Steinke et al. (2023), and indicates the necessity of output set 375 376 optimization for privacy auditing. See Appendix D for more examples of other output distributions.

³⁷⁷

³https://github.com/eth-sri/dp-sniper



Figure 1: Examples where the theoretical optimal output set (purple intervals) for the auditing objective (4) (given m = 100 output samples) contains multiple intervals. This is because the function $\frac{\max\{p(o),q(o)\}}{p(o)+q(o)}$ (red curves), which is indicative of the likelihood ratio, is non-monotonic and contains peaks in the central part of the axis. We experiment for the output distributions p and q in shuffied noisy SGD under (1) Laplace mechanism with $p = 0.5 \cdot \text{Lap}(-1, 1) + 0.5 \cdot \text{Lap}(2, 1)$ and $q = 0.5 \cdot \text{Lap}(-1, 1) + 0.5 \cdot \text{Lap}(2, 1)$ $0.5 \cdot \text{Lap}(-2, 1) + 0.5 \cdot \text{Lap}(1, 1)$, i.e., (11) and 12 with $x_1 = 4, x'_1 = -4, x_2 = x'_2 = 0$ and b = 2; (2) Gaussian mechanism with $p = 0.5 \cdot \mathcal{N}(-2, 1) + 0.5 \cdot \mathcal{N}(1, 1)$ and $q = 0.5 \cdot \mathcal{N}(-1, 1) + 0.5 \cdot \mathcal{N}(0, 1)$, i.e., (9) and 10 with $x_1 = -4$, $x_2 = 0$, $x'_1 = -\frac{4}{3}$, $x'_2 = -\frac{8}{3}$ and $\sigma^2 = \frac{16}{5}$. We also show the best-performing top and bottom output set (green intervals) that is selected according to Steinke et al. (2023), which fails to cover the central regions of the axis that have high values of $\frac{\max\{p(o),q(o)\}}{p(o)+q(o)}$. p(o)+q(o)

401 402 403

404

392

394

397

398 399 400

5.2 Gain of Optimizing the Output Set O

We now experimentally investigate whether choosing the optimal output set enables gain for auditing 405 differential privacy lower bound, compared to prior (suboptimal) strategies (Bichsel et al., 2021; 406 Lu et al., 2022; Steinke et al., 2023) for selecting the output set. Our results are summarized in 407 Figure 2. We first observe that our output set optimization method consistently enables tighter 408 *privacy auditing lower bound* than prior methods (in terms of higher mean and smaller standard 409 deviation) across different specified guarantees ε , as observed in Figure 2 (left). 410

We also observe in Figure 2 (right) that our method is the most advantageous (compared to prior 411 methods) given a small number of auditing samples from the output distributions p and q. This is 412 intuitive because when the number of samples is small, the second stochastic variance term in our 413 optimization objective (4) is dominating. In such scenarios, it is crucial to subselect more auditing 414 samples with the highest absolute log-likelihood ratio statistics, thus making the effect of output set 415 optimization significant. By contrast, prior output set optimization objectives (Bichsel et al., 2021; 416 Lu et al., 2022) do not capture such a sampling variance term, thus leading to suboptimal output sets. 417

Finally, we observe that the amount of auditing gain from the optimal output set is mechanism-418 *dependent*. Specifically, when the number of auditing samples is large enough, all the methods 419 lead to matching privacy lower bounds (to the upper bound $\varepsilon = 1$) in Figure 2b, thus indicating 420 that existing auditing techniques are tight for the mixture of Laplace mechanism. However, for the 421 mixture of Gaussian mechanism (Figure 2d), the auditing gain from our method (compared to other 422 methods) is significant even when the number of auditing samples is large. We hypothesize that this 423 is because mixture of Gaussian distributions incurs lighter tails than mixture of Laplace distributions, 424 thus inducing a higher likelihood ratio in the center of the score domain and necessitating the output 425 set selection beyond the tails. It is an interesting open problem as to what general properties of the 426 mechanisms would make output set optimization more effective in privacy auditing.

427 428

APPLICATION: TIGHTER BLACK-BOX AUDITING OF DP-SGD 6

429 430

In this section, we apply our approximation method (Section 4.3) to compute the optimal output 431 set for auditing black-box differentially private stochastic gradient descent (DP-SGD (Abadi et al.,

463

464

465

467

468 469 470

471

472

473 474



Figure 2: Auditing performance of our method (red) compared with prior methods (Steinke 460 et al., 2023; Bichsel et al., 2021; Lu et al., 2022) (as discussed in Section 5). All evaluation performances are averaged across five random auditing trials. We experiment given $m \in$ 462 $\{100, 200, 400, 800, 1600, 3200, 6400, 12800, 51200\}$ auditing samples from the following output distributions (a) mixture of Laplace mechanism given by (11) and (12) with $x_1 = 4$, $x'_1 = -4$, $x_2 = x'_2 = 0$ and $b \in \{0.25, 0.5, 1, 2, 4, 8\}$ and (b) mixture of Gaussian mechanism given by (9) and (10) with $x_1 = -4$, $x_2 = 0$, $x'_1 = -\frac{4}{3}$, $x'_2 = -\frac{8}{3}$ and $\sigma \in \frac{4}{\sqrt{5}} \cdot \{0.25, 0.5, 1, 2, 4, 8\}$. We 466 compute the DP guarantees ε specified by different noise scale b and σ following (Dwork et al., 2006, Theorem 1) and (Balle & Wang, 2018, Theorem 8) (see Proposition C.1 for the details).

2016)). To reduce the computation cost, we focus on the one-run auditing experiment as introduced by Steinke et al. (2023). To ensure the validity of the auditing lower bound, we use the auditing function in [Steinke et al., 2023, Theorem 5.2] in our output selection Algorithm 1 (Line 6). We show that our method provides a tighter auditing lower bound than existing methods.

Experiment Setting We start with a state-of-the-art DP training setting from (Sander et al., 2023): 475 running DP-SGD on the CIFAR-10 dataset with a fixed privacy budget ($\varepsilon = 8, \delta = 10^{-5}$), with 476 16-4-WideResNet (Zagoruyko, 2016). When given the *full* CIFAR-10 as training dataset, our exper-477 iment reaches 76% test accuracy which matches the state-of-the-art performance reported in Sander 478 et al. (2023); De et al. (2022). Following the setting of black-box input space auditing experiment in 479 Steinke et al. (2023, 6.2), we use loss as the MIA score, and consider the setting where the auditor 480 only has control over the what images to use for training, cannot tweak any intermediate part (e.g., 481 gradient) of the DP-SGD training procedure, and observes only the final trained model. Following 482 Steinke et al. (2023), we experiment on both natural in-distribution data records (actual CIFAR-483 10 records) and canary data records (mislabelled CIFAR-10 records) as training datasets used for auditing. For simplicity, we focus on the setting where all records used for training are used for au-484 diting, i.e., each record is independently included in the training dataset with half probability. This 485 is consistent with the setting of (Steinke et al., 2023, Figure 8).

486 We then perform output set selection by our method and compare it with the techniques in prior 487 works (as discussed in Section 5). To ensure that the comparison is fair, we either used the code 488 released by the authors (for DP-sniper (Bichsel et al., 2021)) or ensured that our implementations 489 of prior baselines enabled matching auditing performance as reported in the prior work (e.g., our 490 observed performance for Steinke et al. (2023) in Figure 3b matches (Steinke et al., 2023, Figure 8) given m = 10000 auditing examples). 491

492 Gain of choosing the optimal output set Figure 3 summarizes our results. We show the auditing 493 lower bounds enabled by output sets \mathcal{O}_{τ} constructed by our method Algorithm 1 as well as the by 494 the heuristic intervals in Steinke et al. (2023). To rule out the effect of estimated intervals overfitting to the empirical Monte Carlo samples from output distributions, for both our method and prior 495 methods, we used a set of fresh samples (that are disjoint from the samples used for estimating the 496 output set) to evaluate the auditing lower bounds. We observe that our method generally provides 497 a better (or comparable) auditing lower bound than the heuristic method in Steinke et al. (2023), in 498 terms of higher mean or smaller variance. This shows that our method is able to approximate the 499 optimal output set that maximizes the privacy loss lower bound, while the heuristic top (bottom) 500 selection method in Steinke et al. (2023) may not be able to achieve this in certain settings. 501

Effect of auditing data examples on the privacy lower bound We observe that the auditing 502 lower bound is higher for the canary dataset (Figure 3b) than the natural dataset (Figure 3a), under 503 the same fixed specified DP guarantee. This is consistent with the intuition that the canary dataset 504 is closer to the worst-case data record which induces high information leakage. This also suggests 505 that the auditing lower bound is sensitive to the dataset used for auditing, and the choice of dataset 506 can significantly affect the auditing performance. We also observe that the gain of our method is 507 higher in the setting with 1000 records (Figure 3b) compared to the setting with 10000 records (Figure 3c). This is consistent with the intuition that the effect of selecting the optimal output set is 508 more significant when the number of MC samples is small, as discussed in Section 5.2. 509



Figure 3: Auditing performance versus specified DP guarantee ε for our method (red) compared prior methods (Steinke et al., 2023; Bichsel et al., 2021; Lu et al., 2022) (as discussed in Section 5). We run DP-SGD on both natural in-distribution actual images and mislabelled canary images in the CIFAR-10 dataset. All performance are averaged across five random auditing trials.

7 CONCLUSION

514

517

518

521 522

523

524

525

526 527

528

529 In this paper, we proposed a framework to compute the optimal output event set that maximizes the 530 privacy lower bound in privacy auditing. We derive optimality guarantee for the output set formed 531 by thresholding likelihood ratio statistics between member and non-member score distributions in 532 the auditing experiment. Through experiments, we show that optimizing the output set effectively 533 tightens the privacy lower bound estimate compared to existing auditing techniques, and provides 534 a more accurate analysis of differentially-private learning algorithms (such as subsampled Laplace 535 and Gaussian mechanisms as well as black-box DP-SGD training). Interestingly, we find that output set optimization is the most effective when given a *restricted number of auditing examples*, or when 537 the likelihood ratio function is non-monotonic or assymetric (e.g., for mechanisms with asymmetric or multi-modal output distributions). Future work includes extending our framework to other per-538 formance metrics of inference attacks, and exploring more accurate approximations of the output set when only given empirical samples from the output distributions.

540 REFERENCES 541

547

551

554

559

560

565

566

567

570

576

580

581

582

- Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and 542 Li Zhang. Deep learning with differential privacy. In Proceedings of the 2016 ACM SIGSAC 543 conference on computer and communications security, pp. 308–318, 2016. 544
- Michael Aerni, Jie Zhang, and Florian Tramèr. Evaluations of machine learning privacy defenses 546 are misleading. arXiv preprint arXiv:2404.17399, 2024.
- Borja Balle and Yu-Xiang Wang. Improving the gaussian mechanism for differential privacy: Ana-548 lytical calibration and optimal denoising. In International Conference on Machine Learning, pp. 549 394-403. PMLR, 2018. 550
- Borja Balle, Giovanni Cherubin, and Jamie Hayes. Reconstructing training data with informed 552 adversaries. In 2022 IEEE Symposium on Security and Privacy (SP), pp. 1138–1156. IEEE, 553 2022.
- Raef Bassily, Adam Smith, and Abhradeep Thakurta. Private empirical risk minimization: Efficient 555 algorithms and tight error bounds. In 2014 IEEE 55th annual symposium on foundations of 556 computer science, pp. 464-473. IEEE, 2014.
- 558 Benjamin Bichsel, Samuel Steffen, Ilija Bogunovic, and Martin Vechev. Dp-sniper: Black-box discovery of differential privacy violations using classifiers. In 2021 IEEE Symposium on Security and Privacy (SP), pp. 391–409. IEEE, 2021.
- 561 Nicholas Carlini, Steve Chien, Milad Nasr, Shuang Song, Andreas Terzis, and Florian Tramer. Mem-562 bership inference attacks from first principles. In 2022 IEEE Symposium on Security and Privacy 563 (SP), pp. 1897–1914. IEEE, 2022.
 - Soham De, Leonard Berrada, Jamie Hayes, Samuel L Smith, and Borja Balle. Unlocking high-accuracy differentially private image classification through scale. arXiv preprint arXiv:2204.13650, 2022.
- 568 Carles Domingo-Enrich and Youssef Mroueh. Auditing differential privacy in high dimensions with 569 the kernel quantum r\'enyi divergence. arXiv preprint arXiv:2205.13941, 2022.
- Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity 571 in private data analysis. In Theory of Cryptography: Third Theory of Cryptography Conference, 572 TCC 2006, New York, NY, USA, March 4-7, 2006. Proceedings 3, pp. 265–284. Springer, 2006. 573
- 574 Cynthia Dwork, Aaron Roth, et al. The algorithmic foundations of differential privacy. Found. 575 Trends Theor. Comput. Sci., 9(3-4):211–407, 2014.
- Chuan Guo, Brian Karrer, Kamalika Chaudhuri, and Laurens van der Maaten. Bounding training 577 data reconstruction in private (deep) learning. In International Conference on Machine Learning, 578 pp. 8056-8071. PMLR, 2022. 579
 - Chuan Guo, Alexandre Sablayrolles, and Maziar Sanjabi. Analyzing privacy leakage in machine learning via multiple hypothesis testing: A lesson from fano. In International Conference on Machine Learning, pp. 11998–12011. PMLR, 2023.
- 583 Matthew Jagielski, Jonathan Ullman, and Alina Oprea. Auditing differentially private machine 584 learning: How private is private sgd? Advances in Neural Information Processing Systems, 33: 585 22205-22216, 2020. 586
 - Bargav Jayaraman and David Evans. Evaluating differentially private machine learning in practice. In 28th USENIX Security Symposium (USENIX Security 19), pp. 1895–1912, 2019.
- 589 Peter Kairouz, Sewoong Oh, and Pramod Viswanath. The composition theorem for differential 590 privacy. In International conference on machine learning, pp. 1376–1385. PMLR, 2015. 591
- William Kong, Andres Munoz Medina, Monica Ribero, and Umar Syed. Dp-auditorium: A large 592 scale library for auditing differential privacy. In 2024 IEEE Symposium on Security and Privacy 593 (SP), pp. 219–219. IEEE Computer Society, 2024.

594 Yiyong Liu, Zhengyu Zhao, Michael Backes, and Yang Zhang. Membership inference attacks by 595 exploiting loss trajectory. In Proceedings of the 2022 ACM SIGSAC Conference on Computer and 596 Communications Security, pp. 2085–2098, 2022. 597 Fred Lu, Joseph Munoz, Maya Fuchs, Tyler LeBlond, Elliott Zaresky-Williams, Edward Raff, Fran-598 cis Ferraro, and Brian Testa. A general framework for auditing differentially private machine learning. Advances in Neural Information Processing Systems, 35:4165-4176, 2022. 600 601 Ilya Mironov. Rényi differential privacy. In 2017 IEEE 30th computer security foundations sympo-602 sium (CSF), pp. 263–275. IEEE, 2017. 603 604 Milad Nasr, Shuang Songi, Abhradeep Thakurta, Nicolas Papernot, and Nicholas Carlin. Adversary 605 instantiation: Lower bounds for differentially private machine learning. In 2021 IEEE Symposium 606 on security and privacy (SP), pp. 866–882. IEEE, 2021. 607 Milad Nasr, Jamie Hayes, Thomas Steinke, Borja Balle, Florian Tramèr, Matthew Jagielski, 608 Nicholas Carlini, and Andreas Terzis. Tight auditing of differentially private machine learning. 609 arXiv preprint arXiv:2302.07956, 2023. 610 611 Jerzy Neyman and Egon Sharpe Pearson. Ix. on the problem of the most efficient tests of statistical 612 hypotheses. Philosophical Transactions of the Royal Society of London. Series A, Containing 613 Papers of a Mathematical or Physical Character, 231(694-706):289–337, 1933. 614 Krishna Pillutla, Galen Andrew, Peter Kairouz, H Brendan McMahan, Alina Oprea, and Sewoong 615 Oh. Unleashing the power of randomization in auditing differentially private ml. arXiv preprint 616 arXiv:2305.18447, 2023. 617 618 Alexandre Sablayrolles, Matthijs Douze, Cordelia Schmid, Yann Ollivier, and Hervé Jégou. White-619 box vs black-box: Bayes optimal strategies for membership inference. In International Confer-620 ence on Machine Learning, pp. 5558–5567. PMLR, 2019. 621 622 Ahmed Salem, Yang Zhang, Mathias Humbert, Pascal Berrang, Mario Fritz, and Michael Backes. Ml-leaks: Model and data independent membership inference attacks and defenses on machine 623 learning models. arXiv preprint arXiv:1806.01246, 2018. 624 625 Tom Sander, Pierre Stock, and Alexandre Sablayrolles. Tan without a burn: Scaling laws of dp-sgd. 626 In International Conference on Machine Learning, pp. 29937–29949. PMLR, 2023. 627 628 Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference at-629 tacks against machine learning models. In 2017 IEEE symposium on security and privacy (SP), pp. 3–18. IEEE, 2017. 630 631 Liwei Song and Prateek Mittal. Systematic evaluation of privacy risks of machine learning models. 632 In 30th USENIX Security Symposium (USENIX Security 21), pp. 2615–2632, 2021. 633 634 Thomas Steinke, Milad Nasr, and Matthew Jagielski. Privacy auditing with one (1) training run. 635 arXiv preprint arXiv:2305.08846, 2023. 636 637 Jasper Tan, Blake Mason, Hamid Javadi, and Richard Baraniuk. Parameters or privacy: A provable 638 tradeoff between overparameterization and membership inference. Advances in Neural Information Processing Systems, 35:17488–17500, 2022. 639 640 Florian Tramer, Andreas Terzis, Thomas Steinke, Shuang Song, Matthew Jagielski, and Nicholas 641 Carlini. Debugging differential privacy: A case study for privacy auditing. arXiv preprint 642 arXiv:2202.12219, 2022. 643 644 Larry Wasserman and Shuheng Zhou. A statistical framework for differential privacy. Journal of 645 the American Statistical Association, 105(489):375-389, 2010. 646 Samuel S Wilks. The large-sample distribution of the likelihood ratio for testing composite hypothe-647

ses. The annals of mathematical statistics, 9(1):60-62, 1938.

- Jiayuan Ye, Aadyaa Maddi, Sasi Kumar Murakonda, Vincent Bindschaedler, and Reza Shokri. Enhanced membership inference attacks against machine learning models. In *Proceedings of the* 2022 ACM SIGSAC Conference on Computer and Communications Security, pp. 3093–3106, 2022.
- Samuel Yeom, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. Privacy risk in machine learning: Analyzing the connection to overfitting. In 2018 IEEE 31st computer security foundations symposium (CSF), pp. 268–282. IEEE, 2018.
- 656 Sergey Zagoruyko. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.
- Santiago Zanella-Beguelin, Lukas Wutschitz, Shruti Tople, Ahmed Salem, Victor Rühle, Andrew Paverd, Mohammad Naseri, Boris Köpf, and Daniel Jones. Bayesian estimation of differential privacy. In *International Conference on Machine Learning*, pp. 40624–40636. PMLR, 2023.

1	Introduction					
2	Related Works Problem Statement					
3						
4	Computing the Optimal Output Set for Privacy Auditing					
	4.1	Optimization Objective for Output Set Selection				
	4.2	Identifying Optimal Output Set for Auditing				
	4.3	Approximating the Optimal Output Set from Empirical Samples				
5	Nun	nerical Experiments: One-Epoch of Shuffled Noisy SGD				
	5.1	Shape of Optimal Output Set				
	5.2	Gain of Optimizing the Output Set \mathcal{O}				
6	Application: Tighter Black-box Auditing of DP-SGD					
7	Conclusion					
/						
A	Deta	ails regarding auditing experiments in the literature				
/ A B	Deta Defe	ails regarding auditing experiments in the literature erred proofs for Section 4				
/ A B	Deta Defe B.1	ails regarding auditing experiments in the literature erred proofs for Section 4 Notations and Results from Prior Works				
/ A B	Deta Defe B.1 B.2	ails regarding auditing experiments in the literature erred proofs for Section 4 Notations and Results from Prior Works				
A B	Deta Defa B.1 B.2 B.3	ails regarding auditing experiments in the literature erred proofs for Section 4 Notations and Results from Prior Works Proof for Proposition 4.1 Proof for Theorem 4.3				
A B C	Deta Defa B.1 B.2 B.3 Defa	ails regarding auditing experiments in the literature erred proofs for Section 4 Notations and Results from Prior Works Proof for Proposition 4.1 Proof for Theorem 4.3 Proofs for Section 5				
A B C D	Deta Defe B.1 B.2 B.3 Defe Mon	ails regarding auditing experiments in the literature erred proofs for Section 4 Notations and Results from Prior Works Proof for Proposition 4.1 Proof for Theorem 4.3 Prevent Proofs for Section 5 re examples for the shape of output set in other mechanisms				
A B C D E	Deta Defa B.1 B.2 B.3 Defa Mon	ails regarding auditing experiments in the literature erred proofs for Section 4 Notations and Results from Prior Works Proof for Proposition 4.1 Proof for Theorem 4.3 Proofs for Section 5 erred Proofs for Section 5 itional Results				
A B C D E	Deta Defa B.1 B.2 B.3 Defa Mon Add E.1	Ails regarding auditing experiments in the literature erred proofs for Section 4 Notations and Results from Prior Works				

In this section, we explain the three components (as described in Section 3) for privacy auditing experiments in more detail.

749Computing Member and Non-member Scores S_+ and S_- Typically, the set of member scores750 S_+ and non-member scores S_- consists of independent Monte Carlo samples Jagielski et al. (2020);751Nasr et al. (2021; 2023); Zanella-Beguelin et al. (2023); Bichsel et al. (2021); Lu et al. (2022) from752the member and non-member score distributions p and q specified by the auditing experiment (i.e.,753via computing the MIA scores on trained model in the auditing experiment and their training and test754data record respectively). An exception is the recent literature of efficient privacy auditing Steinke755et al. (2023); Pillutla et al. (2023), where S_+ and S_- only consists of weakly independent samples756due to the restricted number of trained models (one Steinke et al. (2023) or a few Pillutla et al.

Ref	Inputs	Member scores S_+	Non- member scores S_	Output set O	Auditing function L
Jagielski et al. (2020); Nasr et al. (2021)	Datasets D , D' with k differing (poisoned) data records D_p . Mod- els $\theta_1, \dots, \theta_T \xleftarrow{i.i.d.T} D$, $\theta'_1, \dots, \theta'_T \xleftarrow{i.i.d.T} D'$.	$\begin{array}{r c} Score(\theta_t, x) \\ x \in D_p \\ \text{for } t = \\ 1, \cdots, T \end{array}$	$\begin{array}{l} : \ Score(\theta_t', x_i) \\ x \ \in \ D_p \\ \text{for} \ t \ = \\ 1, \cdots, T \end{array}$	$:(Z,\infty)$ for threshold Z	(Jagielski et al 2020, Algorithr 2)
Steinke et al. (2023)	$\begin{array}{ccc} \text{Datasets} & D_{pop} & \text{and} \\ D_{can}. & \text{Model} \\ \theta & \overleftarrow{\mathcal{T}} & D_{pop} & \cup & D, \\ D & \overleftarrow{\text{i.i.d. inclusion}} & D_{can}. \end{array}$	$Score(\theta, x_i)$ $x_i \in D$	$\begin{array}{ccc} : Score(\theta, x_i) \\ x_i & \in \\ D_{can} \setminus D \end{array}$: Top k_+ and bottom k values of $S_+ \cup S$	(Steinke et al 2023, Theorem 5.2), als restated in Coro lary B and E.2
Ours	-	-	-	Algorithm 1	-

Table 1: Comparison between auditing experiments in the literature and our auditing experiment. Our method only modifies the construction of output set \mathcal{O} , and can be generally applied on top of any existing auditing experiment.

(2023)) in the auditing experiments. In this paper, we will use two standard ways to compute S_+ and S_- in the literature Jagielski et al. (2020); Steinke et al. (2023). Our framework is general and can be applied on top of any way of computing member and non-member scores in existing auditing experiments.

785

802

Determining the output event set \mathcal{O} To obtain higher auditing lower bound, it is crucial to choose 786 an appropriate output set \mathcal{O} for subselecting the member and non-member scores. Intuitively, this 787 is because we would like to avoid inferring membership for outputs that are equally likely to be 788 observed under member distribution p and non-member distribution q. For example, when $p(o) \sim$ 789 $\mathcal{N}(0,1)$ and $q(o) \sim \mathcal{N}(1,1)$ follow Gaussian distributions, the output $o = \frac{1}{2}$ is equally likely to be 790 observed under p and q. By contrast, the output o = 0 is significantly more likely to be observed 791 under p than q. Thus, intuitively we select a subset \mathcal{O} of member scores S_+ and non-member scores 792 S_{-} to increase their distinguishability (i.e., we would like \mathcal{O} to include o = 0 but not $o = \frac{1}{2}$ in the 793 Gaussian distribution example). To this end, prior works Bichsel et al. (2021); Lu et al. (2022) have 794 proposed to train a neural network model to learn the optimal output set \mathcal{O} for subselecting member 795 and non-member scores by maximizing the likelihood difference.

However, such strategies miss an important aspect of privacy auditing experiment: S_+ and S_- are empirical samples (rather than closed-form densities) from the output distribution. Consequently, the lower bound estimate (which is a random variable) for privacy auditing suffers from a significant sampling error when the number of selected samples from S_+ and S_- is small. Thus, to reduce the sampling error, we would like the output set \mathcal{O} to contain as many samples from S_+ and S_- as possible.

Bosigning auditing function *L* To design an auditing function *L* such that Equation (1) is satisfied under $\hat{\varepsilon} := L(S_+, S_-, \mathcal{O}; \delta, \beta)$, prior works generally relies on the constraints of (ε, δ) -differential privacy on the performance of a binary membership inference hypothesis test Wasserman & Zhou (2010); Kairouz et al. (2015). In this paper, we focus on an advantage-based auditing experiment, i.e., the design of the function *L* is based on bounding the advantage of membership inference with the differential privacy guarantee (ε, δ) . Extending our framework to other performance metrics of binary membership inference hypothesis tests such as true positive rate (TPR) and false positive rate (FPR) is an interesting future work.

B DEFERRED PROOFS FOR SECTION 4

812 B.1 NOTATIONS AND RESULTS FROM PRIOR WORKS

We first define the following experiment for distinguishing between two distributions p and q, while performing membership inference on i.i.d. samples from p and q that fall into a fixed output set O.

816 **Definition B.1** (Experiment for distinguishing between distributions p and q). Let S_+ be m i.i.d. 817 Monte Carlo samples from distribution p, and let S_- be m i.i.d. samples from distribution q. Let \mathcal{O} 818 be a fixed output set.

For an arbitrary inference agorithm $\mathcal{M} : \mathbb{R} \to \{-1, +1\}$ that maps an output sample $o \in \mathbb{R}$ to a binray guess in $\{-1, 1\}$ for whether the sample is drawn from p (guesses 1) or q (guesses -1), define the following experiment.

- 1. For $i = 1, \dots, m$
- 823 824 825

0

827

828 829

830

831 832 833

834 835

836

837

839

840

841

842

843

863

i-th example in S₊. Otherwise, the challenger sets o_i to be the *i*-th example in S₋.
(b) Challenger then sends o_i to the adversary.
(c) A begin for a formation of *M*(o_i) if o_i ∉ O

(a) Challenger samples $b_i \xleftarrow{uniform} \{1, -1\}$. If $b_i = 1$, the challenger sets o_i to be the

(c) Adversary performs inference
$$\hat{b}_i = \mathcal{A}(o_i) \coloneqq \begin{cases} \mathcal{M}(o_i) & \text{if } o_i \notin \mathcal{O} \\ 0 & \text{if } o_i \in \mathcal{O} \end{cases} \in \{-1, 0, 1\}.$$

2. Return $(b_i)_{i=1}^m$ and $(\hat{b}_i)_{i=1}^m$.

The above experiment is similar to the membership inference experiment in Yeom et al. (2018, Experiment 1), except for the following differences.

- 1. We used two distributions p and q to astract the distributions for member and non-member scores in the auditing experiment respectively. The randomness of distributions p and qcome from many sources, such as the randomness from the dataset sampling and the output sampling of the DP mechanism. See Table 1 row one for examples of i.i.d. samples from pand q in prototypical auditing experiments (Jagielski et al., 2020; Nasr et al., 2021).
 - 2. We incorporated an output set \mathcal{O} to subselect the output samples for subsequent inference. This changed the guess \hat{b}_i from binary (as in Yeom et al. (2018)) to ternary, where the guess 0 indicates that the output sample o_i is not in the output set \mathcal{O} .

To use the returned values b_i and \hat{b}_i from Definition B.1 to audit approximate DP, we will use the auditing bound from Steinke et al. (2023, Theorem 5.2) that generally holds for inference results in the range of [-1, 1] (rather than standard auditing bounds that only holds for binary guesses, such as Yeom et al. (2018); Jagielski et al. (2020); Nasr et al. (2021)). For the simplicity of the presentation, here we restate and prove a variant of Steinke et al. (2023, Theorem 5.2) that is specialized for the output-set-dependent experiment in Definition B.1 under pure DP and i.i.d. samples.

Lemma B.2 (Variant of Steinke et al. (2023, Proposition) under i.i.d. Monte Carlo samples). Let p and q be two probability distributions over \mathbb{R} that satisfies $e^{-\varepsilon} \leq \frac{p(o)}{q(o)} \leq e^{\varepsilon}$ for any $o \in \mathbb{R}$. Let $\mathcal{O} \subset \mathbb{R}$ be a fixed output set. Let b_i , \hat{b}_i be defined as in experiment Definition B.1. Then conditioned on any fixed value $t \in \{-1, 0, 1\}^m$ for $(\hat{b}_i)_{i=1}^m$, it satisfies that

$$\Pr\left[\sum_{i=1}^{m} \max\{b_i \cdot \hat{b}_i, 0\} \ge v | \hat{b} = t\right] \le \Pr_{S \leftarrow Bernoulli\left(\frac{e^{\varepsilon}}{e^{\varepsilon} + 1}\right)^m} \left[S_i \cdot |t_i| \ge v\right]$$
(13)

Proof. Observe that $b_i \cdot \hat{b}_i$ for $i = 1, \dots, m$ are independent random variables, due to the i.i.d. sampling of S_+ , S_- and $b_i \xleftarrow{uniform} \{-1, 1\}$ across $i = 1, \dots, m$ in Definition B.1. Thus to prove Equation (13), we only need to prove that

$$\frac{1}{1+e^{\varepsilon}} \le \Pr[b_i = 1|\hat{b}_i = t_i] \le \frac{e^{\varepsilon}}{1+e^{\varepsilon}}$$
(14)

=

=

Benote $A_1 = \{o \in \mathbb{R} : \mathcal{A}(o) = 1\}, A_0 = \{o \in \mathbb{R} : \mathcal{A}(o) = 0\}$ and $A_{-1} = \{o \in \mathbb{R} : \mathcal{A}(o) = -1\}$ to be the preimage set of guess 1, 0 and -1 respectively, given membership inference strategy \mathcal{A} . Then by definition, we have that

$$\Pr[b_i = 1|\hat{b}_i = t_i] = \frac{\Pr[b_i = 1, \hat{b}_i = t_i]}{\Pr[b_i = -1, \hat{b}_i = t_i] + \Pr[b_i = 1, \hat{b}_i = t_i]}$$
(15)
$$\Pr[b_i = 1] \cdot \Pr[\hat{b}_i = t_i] + \Pr[\hat{b}_i = t_i]$$

869 870 871

872 873

874 875

876 877 878

879 880

882

883

895 896

899

900

901 902

903

868

$$=\frac{\Pr[b_i = 1] \cdot \Pr[b_i = 1]}{\Pr[b_i = -1] \cdot \Pr[\hat{b}_i = t_i | b_i = -1] + \Pr[b_i = 1] \cdot \Pr[\hat{b}_i = t_i | b_i = 1]}$$
(16)

$$=\frac{0.5 \cdot q(A_{t_i})}{0.5 \cdot p(A_{t_i}) + 0.5 \cdot q(A_{t_i})}$$
(17)

$$\frac{1}{1 + \frac{p(A_{t_i})}{q(A_{t_i})}} \in \left[\frac{1}{1 + e^{\varepsilon}}, \frac{1}{1 + e^{-\varepsilon}}\right]$$
(18)

where the last inequality is due to the assumed condition $e^{-\varepsilon} \leq \frac{p(o)}{q(o)} \leq e^{\varepsilon}$ for any $o \in \mathbb{R}$.

On the one hand, Lemma B.2 can be seen as a simplified variant of Steinke et al. (2023, Proposition 5.1) for auditing i.i.d. Monte Carlo samples. On the other hand, this Lemma generalizes Steinke et al. (2023, Proposition 5.1) from special designs of output set (fixed number of "abstentions") to any fixed choice of output set. This generalization then serves as the basis for proving our auditing lower bound Proposition 4.1 that is specific to the choice of output set.

As an corollary of Lemma B.2, we can obtain the following auditing function for pure DP, under subselected i.i.d. output samples that fall into output set \mathcal{O} .

Corollary B.3. Let $\{b_i\}_{i=1}^m$ and $\{\hat{b}_i\}_{i=1}^m$ be returned by Definition B.1 under output set \mathcal{O} . By applying Lemma B.2 under setting M to be the mechanism that maps $\{b_i\}_{i=1}^m$ to $\{\hat{b}_i\}_{i=1}^m$ in Definition B.1, we obtain following auditing function for approximate DP.

$$L(S_{+}, S_{-}, \mathcal{O}; \delta, \beta) = \max\left\{\varepsilon : \Pr_{\substack{S \leftarrow Bernoulli\left(\frac{e^{\varepsilon}}{e^{\varepsilon}+1}\right)^{m}}}\left[S_{i} \cdot |\hat{b}_{i}| \ge v\right]\right\}$$
(19)

where $v = \sum_{i=1}^{m} \max\{b_i \cdot \hat{b}_i, 0\}.$

Observe that the dependence of the above auditing function (19) on the output set \mathcal{O} is implicit through the dependence of $\{\hat{b}_i\}_{i=1}^m$ on \mathcal{O} in Definition B.1. This makes it hard to understand the optimal choice of output set for maximizing the auditing function $L(S_+, S_-, \mathcal{O}; \delta, \beta)$. By contrast, our paper proves new auditing function in Proposition 4.1 that has explicit dependence on the output set \mathcal{O} , and can be used to analyze the optimal output set for auditing DP.

B.2 PROOF FOR PROPOSITION 4.1

Proof for Proposition 4.1. Let $(b_i)_{i=1}^m$ and $(\hat{b}_i)_{i=1}^m$ be the outputs of experiment Definition B.1 given the *m* i.i.d. samples S_+ from *p* and *m* i.i.d. samples S_- from *q* as inputs, and under the inference algorithm $\mathcal{M} : \mathbb{R} \to \{1, -1\}$ in Definition B.1 defined as follows.

907 908 909

910

$$\mathcal{M}(o) = \begin{cases} 1 & \text{if } o \in A \\ -1 & \text{if } o \in A^c \end{cases}$$
(20)

where $A = \{o \in \mathbb{R} : p(o) \ge q(o)\}$ is the acceptance region for the optimal membership inference attack that maximizes the advantage for distinguishing between p and q (see e.g., (Tan et al., 2022, Proposition 3.1) for the optimality guarantee).

By Lemma B.2, conditioned on any fixed $t \in \{-1, 0, 1\}^m$, the following inequality holds.

916
917
$$\Pr\left[\sum_{i=1}^{m} \max\{b_i \cdot \hat{b}_i, 0\} \ge v | \hat{b} = t\right] \le \Pr_{S \leftarrow \operatorname{Bernoulli}\left(\frac{e^{\varepsilon}}{e^{\varepsilon} + 1}\right)^m} \left[S_i \cdot |t_i| \ge v\right]$$
(21)

Now by applying Hoeffding's inequality on the right-hand-side sequence of bounded random variables $S \leftarrow \text{Bernoulli}\left(\frac{e^{\varepsilon}}{e^{\varepsilon}+1}\right)^m$, we further prove that for any fixed $t \in \{-1, 0, 1\}^m$, it satisfies that

$$\Pr_{\substack{S \leftarrow \text{Bernoulli}\left(\frac{e^{\varepsilon}}{e^{\varepsilon}+1}\right)^{m}} \left[S_{i} \cdot |t_{i}| \geq v\right] \leq \exp\left(-\frac{2\left(v - \frac{e^{\varepsilon}}{e^{\varepsilon}+1} \cdot \sum_{i=1}^{m} |t_{i}|\right)^{2}}{\sum_{i=1}^{m} |t_{i}|}\right)$$
(22)

By setting $v = \frac{e^{\varepsilon}}{e^{\varepsilon}+1} \sum_{i=1}^{m} |t_i| + \ln\left(\frac{1}{\beta}\right)$ in (22), and by plugging the result into (13), we prove that the following inequality holds for any $t \in \{-1, 0, 1\}^m$.

$$\Pr\left[\sum_{i=1}^{m} \max\{b_i \cdot \hat{b}_i, 0\} \ge \frac{e^{\varepsilon}}{e^{\varepsilon} + 1} \sum_{i=1}^{m} |t_i| + \sqrt{\sum_{i=1}^{m} |t_i| \cdot \frac{\ln\left(2/\beta\right)}{2}} \middle| \hat{b} = t \right] \le \frac{\beta}{2}$$
(23)

We now prove another high probability upper bound for $\sum_{i=1}^{m} \max\{b_i \cdot \hat{b}_i, 0\}$ as follows. By the definition of $\hat{b}_i = \mathcal{A}(o_i)$ for \mathcal{A} as defined in (20), and by Bayes rule, we have that for any $i = 1, \dots, m$, the following equality holds.

$$\Pr[b_i \cdot \hat{b}_i = 1 | \hat{b}_i \neq 0] = \frac{\Pr[b_i = 1, S_+(i) \in \mathcal{O} \cap A] + \Pr[b_i = -1, S_-(i) \in \mathcal{O} \cap A^c]}{\Pr[b_i = 1, S_+(i) \in \mathcal{O}] + \Pr[b_i = -1, S_-(i) \in \mathcal{O}]}$$
(24)

$$=\frac{p(\mathcal{O} \cap A) + q(\mathcal{O} \cap A^c)}{p(\mathcal{O}) + q(\mathcal{O})} = \frac{\int_{o \in \mathcal{O}} \max\{p(o), q(o)\}do}{p(\mathcal{O}) + q(\mathcal{O})}$$
(25)

Thus by again using Hoeffding's inequality on the sequence of bounded random variable $b_i \cdot \hat{b}_i$ for $i = 1, \dots, m$, we have that

$$\Pr\left[\sum_{i=1}^{m} \max\{b_i \cdot \hat{b}_i, 0\} \le \frac{\int_{o \in \mathcal{O}} \max\{p(o), q(o)\} do}{p(\mathcal{O}) + q(\mathcal{O})} \cdot \sum_{i=1}^{m} |t_i| - \sqrt{\sum_{i=1}^{m} |t_i| \cdot \frac{\ln\left(2/\beta\right)}{2}} \middle| \hat{b} = t \right] \le \frac{\beta}{2}$$

$$(26)$$

By combining (23) and (26) using union bound, we have that

$$\Pr\left[\frac{\int_{o\in\mathcal{O}}\max\{p(o),q(o)\}do}{p(\mathcal{O})+q(\mathcal{O})}\cdot\sum_{i=1}^{m}|t_{i}|-\sqrt{\sum_{i=1}^{m}|t_{i}|\cdot\frac{\ln\left(2/\beta\right)}{2}}\geq\frac{e^{\varepsilon}}{e^{\varepsilon}+1}\sum_{i=1}^{m}|t_{i}|+\sqrt{\sum_{i=1}^{m}|t_{i}|\cdot\frac{\ln\left(2/\beta\right)}{2}}\Big|\hat{b}=t\right]\leq\beta$$

By rearranging the terms in the above inequality, we have that

$$\Pr\left[\varepsilon \le \phi\left(\frac{\int_{o\in\mathcal{O}} \max\{p(o), q(o)\}do}{p(\mathcal{O}) + q(\mathcal{O})} - \sqrt{\frac{2\ln(2/\beta)}{\sum_{i=1}^{m} |t_i|}}\right) \left| \hat{b} = t \right] \le \beta$$
(27)

where $\phi(y) = \log\left(\frac{y}{1-y}\right)$ for $y \in (0,1)$ is the logit function. By observing that $\hat{b}_i = \mathcal{A}(o_i) \neq 0$ if any only if $o_i \in \mathcal{O}$, we have that $\sum_{i=1}^m |t_i| \leq |S_+ \cap \mathcal{O}| + |S_- \cap \mathcal{O}|$. By plugging this into (27) and taking maximum over all possible $t \in \{-1,0,1\}^m$, we obtain the bound in the statement (3). \Box *Remark* B.4. Proposition 4.1 generalizes prior advantage-based auditing functions Yeom et al.

(2018); Steinke et al. (2023) from specific designs of output set to any fixed choice of output set. For example, the advantage-based auditing function (Yeom et al., 2018, Theorem 1) is equivalent (up to monotonic scaling) to Equation (3) under setting the whole output domain as output set, i.e., $\mathcal{O} = \mathbb{R}$. Similarly, by setting the output set $\mathcal{O} = (-\infty, o_{k_-}) \cup (o_{k_+}, +\infty)$ where o_{k_-} and o_{k_+} are the bottom- k_- score and top- k_+ score in $S_- \cup S_+$, we recover the auditing function used under the abstention strategy in Steinke et al. (2023, Algorithm 1, Proposition 5.1).

Remark B.5. Proposition 4.1 proves an auditing function that explicitly depends on the output set
 O, rather than implicitly depending on the output set as in prior works Steinke et al. (2023) (as discussed after Lemma B.2). This then allows us to analyze the optimal output set for auditing DP.

972 B.3 PROOF FOR THEOREM 4.3

Proof for Theorem 4.3. We first prove the existence of \mathcal{O}_{τ} such that $p(\mathcal{O}_{\tau}) + q(\mathcal{O}_{\tau}) = p(\mathcal{O}) + q(\mathcal{O})$, 975 for any feasible output set \mathcal{O} .

977 By observing that p and q are continuous output distributions on \mathbb{R} , we prove that the following 978 function is continuous with respect to $\tau \ge 0$.

$$E(\tau) = p(\mathcal{O}_{\tau}) + q(\mathcal{O}_{\tau}) \text{ where } \mathcal{O}_{\tau} = \{x \in \mathbb{R} : \left|\log\frac{p(x)}{q(x)}\right| \ge \tau\}$$

By definition, we further have that E(0) = 2, $\lim_{\tau \to +\infty} E(\tau) = 0$, and $p(\mathcal{O}) + q(\mathcal{O}) \in [0, 2]$. Thus by using intermediate value theorem for continuous function, we prove that for any feasible output set \mathcal{O} , there exists $\tau \in \mathbb{R}$, such that $p(\mathcal{O}) + q(\mathcal{O}) = p(\mathcal{O}_{\tau}) + q(\mathcal{O}_{\tau})$.

985 We now prove the optimality guarantee. Let \mathcal{O} be any output set that satisfies $p(\mathcal{O}) + q(\mathcal{O}) = p(\mathcal{O}_{\tau}) + q(\mathcal{O}_{\tau})$. By definition of \mathcal{O}_{τ} in Definition 4.2, we have that

$$\left(\mathbf{1}_{x\in\mathcal{O}_{\tau}}-\mathbf{1}_{x\in\mathcal{O}}\right)\cdot\left(\max\{p(x),q(x)\}-\frac{e^{\tau}}{1+e^{\tau}}\cdot\left(p(x)+q(x)\right)\right)\geq0$$
(28)

holds for any $x \in \mathbb{R} \subset A$ where A is a ignorable set.

Doing the integration over $x \in \mathbb{R}$, we immediately have that

$$\int_{x \in \mathcal{O}_{\tau}} \max\{p(x), q(x)\} dx - \frac{e^{\tau}}{1 + e^{\tau}} \cdot (p(\mathcal{O}_{\tau}) + q(\mathcal{O}_{\tau}))$$
(29)

$$\geq \int_{x \in \mathcal{O}} \max\{p(x), q(x)\} dx - \frac{e^{\tau}}{1 + e^{\tau}} \cdot (p(\mathcal{O}) + q(\mathcal{O}))$$
(30)

By plugging the condition that $p(\mathcal{O}) + q(\mathcal{O}) = p(\mathcal{O}_{\tau}) + q(\mathcal{O}_{\tau})$ into the constraints, we have that

$$\int_{x \in \mathcal{O}_{\tau}} \max\{p(x), q(x)\} dx \ge \int_{x \in \mathcal{O}} \max\{p(x), q(x)\} dx \tag{31}$$

Thus by definition of auditing bound for output set in Equation (4), we have that $\hat{\varepsilon}(\mathcal{O}_{\tau}; p, q) \geq \hat{\varepsilon}(\mathcal{O}; p, q)$. \Box

Remark B.6 (Similarity in proof technique to Neyman-Pearson Lemma). The proof technique, which constructs an indicator function that is always non-negative (eq 21), and then performs integration (eq 22 and 23), is the standard technique used for poving Neyman-Pearson Lemma.

C DEFERRED PROOFS FOR SECTION 5

Proposition C.1. Let $b_1, b_2, \sigma_1, \sigma_2 \ge 0$, $\theta_0 \in \mathbb{R}$ and $\eta \le 1$ be fixed. Assume that the input dataset D has bounded domain, in that there exists $r \ge 0$ such that $|x| \le r$ for any $x \in D$. Define division over zero as infinity, i.e., $\frac{1}{0} = \infty$. Then

• If $Z_1 \sim Lap(0, b_1)$ and $Z_2 \sim Lap(0, b_2)$, then the mechanism with output distribution (8) satisfies ε -DP with $\varepsilon = \frac{r\eta(2-\eta)}{\max\{\eta(1-\eta)b_1, \eta b_2\}}$.

• If $Z_1 \sim \mathcal{N}(0, \sigma_1^2)$ and $Z_2 \sim \mathcal{N}(0, \sigma_2^2)$, then the mechanism with output distribution (8) satisfies (ε, δ) -DP with $\delta = \bar{\Phi}\left(\frac{r(2-\eta)}{2\sigma} - \frac{\varepsilon\sigma}{r(2-\eta)}\right) - e^{\varepsilon} \cdot \bar{\Phi}\left(-\frac{r(2-\eta)}{2\sigma} - \frac{\varepsilon\sigma}{r(2-\eta)}\right)$ for any $\varepsilon \geq 0$, where $\sigma = \sqrt{(1-\eta)^2\sigma_1^2 + \sigma_2^2}$ and Φ denotes the cumulative distribution function of the standard normal distribution.

Proof. It suffices to apply (Dwork et al., 2006, Theorem 1) and (Balle & Wang, 2018, Theorem 8) after observing that the ℓ_1 and ℓ_2 sensitivity of the update in (8) is $\max_{x_{s_1}, x_{s_2}} |-\eta \cdot (1-\eta)x_{s_1} + \eta x_{s_2}| \le \eta(2-\eta) \cdot r$ for $\eta \le 1$.



Figure 4: Examples where our selected optimal output set (purple intervals) is asymmetric and 1043 consists of multiple intervals, given a small number of m = 100 auditing MCMC samples. We 1044 show the function $\frac{\max\{p(o),q(o)\}}{p(o)+q(o)}$ (that is indicative of the absolute value of log-likelihood ratio) p(o)+q(o)1045 for the output distributions p and q in shuffled noisy SGD under (1) Cauchy distributions with 1046 p = Cauchy(-1, 0.1) and q = Cauchy(1, 3); (2) F distributions with p = F(3, 10) and q = F(3, 10)1047 F(5, 10). We also show the best-performing top and bottom output set (green intervals) that is 1048 selected according to Steinke et al. (2023), which fails to cover the central regions of the axis that have high $\frac{\max\{p(\bar{x}),q(x)\}}{\max\{p(\bar{x}),q(x)\}}$ 1049 p(x)+q(x)1050



1069 Figure 5: Examples where our selected optimal output set (purple intervals) is asymmetric and 1070 consists of multiple intervals, given a small number of m = 100 auditing MCMC samples. We show the function $\frac{\max\{p(o), q(o)\}}{p(o) + q(o)}$ (that is indicative of the absolute value of log-likelihood ratio) 1071 p(o)+q(o)1072 for the output distributions p and q in shuffled noisy SGD under (1) Logistic distributions with 1073 p = Logistic(-1, 0.1) and q = Logistic(1, 3); (2) Chisquared distributions with $p = \chi_1^2(3)$ and 1074 $q = \chi_2^2(1)$, i.e., p follows a non-central Chi-squared distribution with 1 degree of freedom and non-1075 centrality parameter 3, while q has 2 degrees of freedom and non-centrality parameter 1. We also show the best-performing top and bottom output set (green intervals) that is selected according to Steinke et al. (2023), which fails to cover the central regions of the axis that have high $\frac{\max\{p(x),q(x)\}}{p(x)+p(x)}$. 1077 p(x)+q(x)1078

1080
1081
1082DMORE EXAMPLES FOR THE SHAPE OF OUTPUT SET IN OTHER
MECHANISMS

E ADDITIONAL RESULTS

1086 E.1 APPLYING ALGORITHM 1 TO AUDIT APPROXIMATE DIFFERENTIAL PRIVACY

In this section, we give examples of applying our output set selection Algorithm 1 to audit approximate (ε, δ) -differential privacy. To ensure that the auditing lower bound is valid for approximate DP, in Line 6 of Algorithm 1, we need to use the following bound from Steinke et al. (2023, Theorem 5.2) that is valid for auditing approximate DP using subselected output scores.

Theorem E.1. Steinke et al. (2023, Theorem 5.2) Let $M : \{-1, +1\}^m \rightarrow [-1, +1]^m$ satisfy (ε, δ) -DP. Let $T = M(S) \in [-1, +1]^m$. Then, for any $v \in \mathbb{R}$,

1094 1095

1098

1099

1100 1101 1102

1084

$$\Pr_{S \leftarrow \{-1,+1\}^m, T \leftarrow M(S)} \left[\sum_{i=1}^m \max\{0, T_i \cdot S_i\} \ge v \right] \le \beta(\varepsilon) + \alpha(\varepsilon) \cdot 2m \cdot \delta \tag{32}$$

where

$$\beta(\varepsilon) = \Pr_{\tilde{W}^*} \left[\tilde{W}^* \ge v \right],\tag{33}$$

$$\alpha(\varepsilon) = \max\left\{\frac{1}{i} \left(\Pr_{\tilde{W}^*}[\check{W}^* \ge v - i] - \beta(\varepsilon) : i \in \{1, 2, \cdots, m\}\right)\right\}.$$
(34)

Here \check{W}^* is any distribution on \mathbb{R} that stochastically dominates $\check{W}(t) \coloneqq \sum_{i=1}^m \check{S}_i |t_i|$ for $\check{S} \leftarrow$ Bernoulli $\left(\frac{e^{\varepsilon}}{e^{\varepsilon}+1}\right)^m$ for all t in the support of T.

By applying Theorem E.1 to Definition B.1, we can obtain the following auditing function for approximate DP under subselected scores from output distributions by output set O.

Corollary E.2 (Auditing function for approximate DP). Let $\{b_i\}_{i=1}^m$ and $\{\hat{b}_i\}_{i=1}^m$ be as defined in Definition B.1 under output set \mathcal{O}_{τ} . By applying Theorem 4.3 under setting M to be the mechanism that maps $\{b_i\}_{i=1}^m$ to $\{\hat{b}_i\}_{i=1}^m$, we obtain following auditing function for approximate DP.

$$L(S_{+}, S_{-}, \mathcal{O}_{\tau}; \delta, \beta) = \max \left\{ \varepsilon : \beta(\varepsilon) + \alpha(\varepsilon) \cdot 2m \cdot \delta \le \beta \right\}$$
(35)

where $\beta(\varepsilon)$ and $\alpha(\varepsilon)$ are defined in Theorem E.1 under setting $v = \sum_{i=1}^{m} \max\{b_i \cdot \hat{b}_i, 0\}$.

Auditing Performance for Approximate DP Our results are summarized in Figure 6. We observe that our method outperforms the prior methods (Steinke et al., 2023) and the whole domain baseline in terms of auditing performance. We also observe that the audited lower bound $\hat{\varepsilon}$ monotonically decreases as δ increases, which is consistent with the theoretical guarantee (dashed line in Figure 6).

1121

1112 1113

1115

- 1122
- 1123
- 1124 1125
- 1126

1127

- 1128
- 1129
- 1130
- 1131

1132

