

AUDITING f -DIFFERENTIAL PRIVACY IN ONE RUN

Anonymous authors

Paper under double-blind review

ABSTRACT

Empirical auditing has emerged as a means of catching some of the flaws in the implementation of privacy-preserving algorithms. Existing auditing mechanisms, however, are either computationally inefficient – requiring multiple runs of the machine learning algorithms — or suboptimal in calculating an empirical privacy. In this work, we present a tight and efficient auditing procedure and analysis that can effectively assess the privacy of mechanisms. Our approach is efficient; similar to the recent work of Steinke, Nasr, and Jagielski (2023), our auditing procedure leverages the randomness of examples in the input dataset and requires only a single (training) run of the target mechanism. And it is more accurate; we provide a novel analysis that enables us to achieve tight empirical privacy estimates by using the hypothesized f -DP curve of the mechanism, which provides a more accurate measure of privacy than the traditional ϵ, δ differential privacy parameters. We use our auditing procedure and analysis to obtain empirical privacy, demonstrating that our auditing procedure delivers tighter privacy estimates.

1 INTRODUCTION

Differentially private machine learning (Chaudhuri et al., 2011; Abadi et al., 2016) has emerged as a principled solution to learning models from private data while still preserving privacy. Differential privacy (Dwork, 2006) is a cryptographically motivated definition, which requires an algorithm to possess certain properties: specifically, a randomized mechanism is differentially private if it guarantees that the participation of any single person in the dataset does not impact the probability of any outcome by much.

Enforcing this guarantee requires the algorithm to be carefully designed and rigorously analyzed. The process of designing and analyzing such algorithms is prone to errors and imperfections as has been noted in the literature Tramer et al. (2022). A result of this is that differentially private mechanisms may not perform as intended, either offering less privacy than expected due to flaws in mathematical analysis or implementation, or potentially providing stronger privacy guarantees that are not evident through a loose analysis.

Empirical privacy auditing (Ding et al., 2018; Nasr et al., 2023; Jagielski et al., 2020) has emerged as a critical tool to bridge this gap. By experimentally assessing the privacy of mechanisms, empirical auditing allows for the verification of privacy parameters. Specifically, an audit procedure is a randomized algorithm that takes an implementation of a mechanism M , runs it in a black-box manner, and attempts to test a privacy hypothesis (such as, a differential privacy parameter). The procedure outputs 0 if there is sufficient evidence that the mechanism does not satisfy the hypothesized guarantees and 1 otherwise. The audit mechanism must possess two essential properties: 1) it must have a *provably* small false-negative rate, ensuring that it would not erroneously reject a truly differentially private mechanism, with high probability; 2) it needs to *empirically* exhibit a “reasonable” false positive rate, meaning that when applied to a non-differentially private mechanism, it would frequently reject the privacy hypothesis. The theoretical proof of the false positive rate is essentially equivalent to privacy accounting (Abadi et al., 2016; Dong et al., 2019; Mironov, 2017), which is generally thought to be impossible in a black-box manner (Zhu et al., 2022).

The prior literature on empirical audits of privacy consists of two lines of work, each with its own set of limitations. The first line of work (Ding et al., 2018; Jagielski et al., 2020; Tramer et al., 2022; Nasr et al., 2023) runs a differentially private algorithm multiple times to determine if the privacy

054 guarantees are violated. This is highly computationally inefficient for most private machine learning
 055 use-cases, where running the algorithm a single time involves training a large model.

056
 057 Recent work (Steinke et al., 2023) remove this limitation by proposing an elegant auditing method that
 058 runs a differentially private training algorithm a single time. In particular, they rely on the randomness
 059 of training data to obtain bounds on the false negative rates of the audit procedure. A key limitation
 060 of the approach in Steinke et al. (2023) is that their audit procedure is sub-optimal in the sense that
 061 there is a relatively large gap between the true privacy parameters of mainstream privacy-preserving
 062 algorithms (e.g., Gaussian mechanism) and those reported by their auditing algorithm.

063 In this work, we propose a novel auditing procedure that is computationally efficient and accurate. Our
 064 method requires only a single run of the privacy mechanism¹ and leverages the f -DP curve (Dong
 065 et al., 2019), which allows for a more fine-grained accounting of privacy than the traditional reliance
 066 on ϵ, δ parameters. By doing so, we provide a tighter empirical assessment of privacy.

067 We experiment with our approach on both simple Gaussian mechanisms as well as a model trained
 068 on real data with DP-SGD. Our experiments show that our auditing procedure can significantly
 069 outperform that of Steinke et al. (2023) (see Figure 1). This implies that better analysis may enable
 070 relatively tight auditing of differentially privacy guarantees in a computationally efficient manner in
 071 the context of large model training.

072 **Technical overview:** We briefly summarize the key technical components of our work and compare
 073 it with that of Steinke et al. (2023). Their auditing procedure employed a game similar to a
 074 membership inference process: the auditor selects a set of canaries and, for each canary, decides
 075 whether to inject it into the training set with independent probability 0.5. Once model training is
 076 completed, the auditor performs a membership inference attack to determine whether each canary
 077 was included. The number of correct guesses made by the adversary in this setting forms a random
 078 variable. The key technical contribution of Steinke et al. was to establish a tail bound on this random
 079 variable for mechanisms satisfying (ϵ) -DP. Specifically, they demonstrated that the tail of this random
 080 variable is bounded by that of a binomial distribution, $\text{binomial}(n, p)$, where n is the number of
 081 canaries and $p = \frac{e^\epsilon}{e^\epsilon + 1}$. To extend this analysis to approximate DP mechanisms, they further showed
 082 that the probability of the adversary’s success exceeding this tail bound is at most $O(n \cdot \delta)$.

083 Steinke et al. highlighted a limitation in their approach in auditing specific mechanisms, such as
 084 the Gaussian mechanism. They correctly argue that simplifying the mechanism’s behavior to just
 085 two parameters, (ϵ, δ) , results in sub-optimal auditing of specific mechanisms. In other words, the
 086 effectiveness of membership inference attacks against the Gaussian mechanism differs significantly
 087 from predictions based solely on the (ϵ, δ) parameters. To overcome this limitation, we propose
 088 auditing the entire privacy curve of a mechanism, rather than focusing solely on (ϵ, δ) . Our solution
 089 involves three key technical steps:

- 091 1. We derive an upper bound on the adversary’s success in correctly guessing a specific canary
 092 for mechanisms satisfying f -DP. This bound is an improved version of the result by Hayes
 093 et al. (2023) for bounding training data reconstruction in DP mechanisms. However, this is
 094 insufficient, as the adversary’s guesses could be dependent, potentially leading to correlated
 095 successes (e.g., correctly or incorrectly guessing all samples).
- 096 2. To address the issue of dependency, we refine our analysis by defining p_i as the probability
 097 of the adversary making exactly i correct guesses. We derive a recursive relation that bounds
 098 p_i based on p_1, \dots, p_{i-1} . This recursive bound is the main technical novelty of our work.
 099 To derive this bound, we consider two conditions: the adversary correctly guesses the first
 100 canary or not. In the first case, we use our analysis from Step 1 to bound the probability of
 101 making $i - 1$ correct guesses given that the first guess was correct. For the incorrect guess
 102 case, we perform a combinatorial analysis to eliminate the condition. This analysis uses
 103 the fact that shuffling of the canaries does not change the probabilities of making i correct
 104 guesses. We note that it is crucial not to use the analysis of Step 1 for both cases. This
 105 is because the analysis of Step 1 cannot be tight for both cases at the same time. Finally,

106 ¹In the context of privacy-preserving training of machine learning models, the privacy mechanism refers
 107 to the training algorithm. Therefore, when we mention a single run, we are specifically referring to a single
 execution of the training algorithm, not the inference algorithm.

leveraging the convexity of trade-off functions and applying Jensen’s inequality, we derive our final recursive relation. To the best of our knowledge, This combination of trade-off function with shuffling is a new technique and could have broader applications.

3. Finally, we design an algorithm that takes advantage of the recursive relation to numerically calculate an upper bound on the tail of the distribution. The algorithm is designed carefully so that we do not need to invoke the result of step 2 for very small events.

We also generalize our analysis to a broader notion of canary injection and membership inference. Specifically, we utilize a reconstruction game where the auditor can choose among k options for each canary point, introducing greater entropy for each choice. This generalization allows for auditing mechanisms with fewer canaries.

In the rest of the paper, we first introduce the notions of f -DP and explain what auditing based on f -DP entails. We then present our two auditing procedures, which are based on membership inference and reconstruction attacks (Section 2). In Section 3, we provide a tight analysis of our audit’s accuracy based on f -DP curves. Finally, in Section 4, we describe the experimental setup used to compare the bounds.

2 AUDITING f - DIFFERENTIAL PRIVACY

Auditing privacy involves testing a ”privacy hypothesis” about an algorithm M . Different mathematical forms can be used for a ”privacy hypothesis,” but they all share the common characteristic of being about an algorithm/mechanism M . For example, one possible hypothesis is that applying SGD with specific hyperparameters satisfies some notion of privacy. With this in mind, the privacy hypothesis are often mathematical constraints on the sensitivity of the algorithm’s output to small changes in its input. The most well-known definition among these is (approximate) differential privacy.

Definition 1. A mechanism M is (ϵ, δ) -DP if for all neighboring datasets $\mathcal{S}, \mathcal{S}'$ with $|\mathcal{S} \Delta \mathcal{S}'| = 1$ and all measurable sets T , we have $\Pr[M(\mathcal{S}) \in T] \leq e^\epsilon \Pr[M(\mathcal{S}') \in T] + \delta$.

In essence, differential privacy ensures that the output distribution of the algorithm does not heavily depend on a single data point. Based on this definition, one can hypothesize that a particular algorithm satisfies differential privacy with certain ϵ and δ parameters. Consequently, auditing differential privacy involves designing a test for this hypothesis. We will later explore the desired properties of such an auditing procedure. However, at present, we recall a stronger notion of privacy known as f -differential privacy.

Notation For a function $f: X \rightarrow \mathbb{R}$ we use \bar{f} to denote the function $\bar{f}(x) = 1 - f(x)$.

Definition 2. A mechanism \mathcal{M} is f -DP if for all neighboring datasets $\mathcal{S}, \mathcal{S}'$ and all $|\mathcal{S} \Delta \mathcal{S}'| = 1$ measurable sets T we have

$$\Pr[M(\mathcal{S}) \in T] \leq \bar{f}(\Pr[M(\mathcal{S}') \in T]).$$

Note that this definition generalizes the notion of approximate differential privacy by allowing a more complex relation between the probability distributions of $M(\mathcal{S})$ and $M(\mathcal{S}')$. The following proposition shows how one can express approximate DP as an instantiation of f -DP.

Proposition 3. A mechanism is (ϵ, δ) -DP if it is f -DP with respect to $\bar{f}(x) = e^\epsilon \cdot x + \delta$.

Although the function f could be an arbitrary function, without loss of generality, we only consider a specific class of functions in this notion.

Remark 4. Whenever we say that a mechanism satisfies f -DP, we implicitly imply that f is a valid trade-off function. That is, f is defined on domain $[0, 1]$ and has a range of $[0, 1]$. Moreover, f is a decreasing and convex with $f(x) \leq 1 - x$ for all $x \in [0, 1]$. We emphasize that this is without loss of generality. That is, if a mechanism is f -DP for an arbitrary function $f: [0, 1] \rightarrow [0, 1]$, then it is also f' -DP for valid trade-off function f' with $f'(x) \leq f(x)$ for all $x \in [0, 1]$ (See Proposition 2.2 in Dong et al. (2019)).

Definition 5 (Order of f -DP curves). For two trade-off functions f_1 and f_2 , we say f_1 is more private than f_2 and denote it by $f_1 \geq f_2$ iff $f_1(x) \geq f_2(x)$ for all $x \in [0, 1]$. Also, for a family of trade-off

functions F , we use $\text{maximal}(F)$ to denote the set of maximal elements w.r.t to the privacy relation. Note that F could be a partial ordered set, and the set of maximal points could have more than a single element.

Now that we have defined our privacy hypothesis, we can turn our attention to auditing these notions.

Definition 6 (Auditing f -DP). *An audit procedure takes the description of a mechanism \mathcal{M} , a trade-off function f , and outputs a bit that determines whether the mechanism satisfies f -DP or not. We define the audit procedure as a two-step procedure.*

- *game*: $M \rightarrow O$, In this step, the auditor runs a potentially randomized experiment/game using the description of mechanism $\mathcal{M} \in M$ and obtains some observation $o \in O$.
- *evaluate* : $O \times F \rightarrow \{0, 1\}$, In this step, the auditor will output a bit b based on an observation o and a trade-off function f . This audit operation tries to infer whether the observation o is “likely” for a mechanism that satisfies f -DP.

The audit procedure is ψ -accurate if for all mechanism \mathcal{M} that satisfy f -DP, we have

$$\Pr_{o \leftarrow \text{game}(\mathcal{M})} [\text{evaluate}(o, f) = 1] \geq \psi.$$

Note that we are defining the accuracy only for positive cases. This is the only guarantee we can get from running attacks. For guarantees in negative cases, we need to perform a proper accounting of the mechanism (Wang et al., 2023).

Auditing f -DP vs DP: f -DP can be viewed as a collection of DP parameters, where instead of considering (ϵ, δ) as fixed scalars, we treat ϵ as a function of δ . For any $\delta \in [0, 1]$, there exists an $\epsilon(\delta)$ such that the mechanism satisfies $(\epsilon(\delta), \delta)$ -DP. The f -DP curve effectively represents the entire privacy curve rather than a single (ϵ, δ) pair. Thus, auditing f -DP can be expected to be more effective, as there are more constraints that need to be satisfied. A naive approach for auditing f -DP is to perform an audit for approximate DP at each (ϵ, δ) value along the privacy curve, rejecting if any of the audits fail. However, this leads to suboptimal auditing performance. First, the auditing analysis involves several inequalities that bound the probabilities of various events using differential privacy guarantees. The probability of these events could take any number between $[0, 1]$. Using a single (ϵ, δ) value to bound the probability of all these events cannot be tight because the linear approximation of privacy curve is tight in at most a single point. Hence, the guarantees of (ϵ, δ) -DP cannot be simultaneously tight for all events. However, with f -DP, we can obtain tight bounds on the probabilities of all events simultaneously. Second, For each (ϵ, δ) we have a small possibility of incorrectly rejecting the privacy hypothesis. So if we audit privacy for $(\epsilon(\delta), \delta)$ independently, we will reject any privacy hypothesis with probability 1.0. This challenge can be potentially resolved by using correlated randomness, but that requires a new analysis.

Next, we formally define the notion of empirical privacy (Nasr et al., 2021) based on an auditing procedure. This notion essentially provides the best privacy guarantee that is not violated by auditors’ observation from a game setup.

Definition 7 (Empirical Privacy). *Let $(\text{game}, \text{evaluate})$ be an audit procedure. We define the empirical privacy random variable for a mechanism \mathcal{M} , w.r.t a family F of trade-off functions, to be the output of the following process. We first run the game to obtain observation $o = \text{game}(\mathcal{M})$. We then construct*

$$F_o = \text{maximal}(\{f \in F; \text{evaluate}(o, f) = 1\})$$

where the maximal set is constructed according to Definition 5. Then, the empirical privacy of the mechanism at a particular δ is defined as

$$\epsilon(\delta) = \min_{f \in F_o} \max_{x \in [0, 1]} \frac{1 - f(x) - \delta}{x}.$$

Note that the empirical privacy $\epsilon(\delta)$ is a function of the observation o . Since, o itself is a random variable, then $\epsilon(\delta)$ is also a random variable.

How to choose the family of trade-off functions? The family of trade-off functions should be chosen based on the expectations of the true privacy curve. For example, if one expects the privacy curve of a mechanism to be similar to that of a Gaussian mechanism, then they would choose the set of all trade-off functions imposed by a Gaussian mechanism as the family. For example, many believe that in the hidden state model of privacy (Ye & Shokri, 2022), the final model would behave like a Gaussian mechanism with higher noise than what is expected from the accounting in the white-box model (where we assume we release all the intermediate models). Although we may not be able to prove this hypothesis, we can use our framework to calculate the empirical privacy, while assuming that the behavior of the final model would be similar to that of a Gaussian mechanism.

2.1 GUESSING GAMES

Here, we introduce the notion of guessing games which is a generalization of membership inference attacks (Nasr et al., 2023), and closely resembles the reconstruction setting introduced in Hayes et al. (2023).

Definition 8. Consider a mechanism $M : [k]^m \rightarrow \Theta$. In a guessing game we first sample an input dataset $\mathbf{u} \in [k]^m$ from an arbitrary distribution. We run the mechanism to get $\theta \sim M(\mathbf{u})$. Then a guessing adversary $A : \Theta \rightarrow ([k] \cup \{\perp\})^m$ tries to guess the input to the mechanism from the output. We define

- the number of guesses by $c' = \sum_{i=1}^m \mathbf{I}(A(\theta)_i \neq \perp)$
- and the number of correct guesses by $c = \sum_{i=1}^m \mathbf{I}(A(\theta)_i = \mathbf{u}_i)$.

Then we output (c, c') as the output of the game.

These guessing games are integral to our auditing strategies. We outline two specific ways to instantiate the guessing game. The first procedure is identical to that described in the work of Steinke et al. (2023) and resembles membership inference attacks. The second auditing algorithm is based on the reconstruction approach introduced by Hayes et al. (2023). In Section 3, we present all of our results in the context of the general notion of guessing games, ensuring that our findings extend to both the membership inference and reconstruction settings.

Auditing by membership inference: Algorithm 1 describes a game setup based on membership inference attacks. In this setup, we have a fixed training set \mathcal{T} and a set of canaries \mathcal{C} . We first sample a subset \mathcal{S} of the canaries using poisson sampling. Then we run the mechanism \mathcal{M} on $\mathcal{T} \cup \mathcal{S}$ to get a model $\theta \sim \mathcal{M}(\mathcal{T} \cup \mathcal{S})$. Then the adversary A inspects θ and tries to find examples that were present in \mathcal{S} . Observe that this procedure is a guessing game with $k = 2$ and $m = |\mathcal{C}|$. This is simply because the adversary is guessing between two choices for each canary, it is either included or not included. Note that this procedure is modular, we can use any \mathcal{T} and \mathcal{C} for the training set and canary set. We can also use any attack algorithm A .

We note that membership inference attacks have received a lot of attention recently (Homer et al., 2008; Shokri et al., 2017; Leino & Fredrikson, 2020; Bertran et al., 2024; Hu et al., 2022; Matthew et al., 2023; Duan et al., 2024; Zarifzadeh et al., 2023). These attack had a key difference from our attack setup and that is the fact that there is only a single example that the adversary is trying to make the inference for. Starting from the work of (Shokri et al., 2017), researchers have tried to improve attacks in various settings (Ye et al., 2022; Zarifzadeh et al., 2023). For example, using calibration techniques has been an effective way to improve membership inference attacks (Watson et al., 2021; Carlini et al., 2022). Researchers have also changed their focus from average case performance of the attack to the tails of the distribution and measured the precision at low recall values (Ye et al., 2022; Nasr et al., 2021).

A substantial body of research has also explored the relationship between membership inference attacks and differential privacy (Sablayrolles et al., 2019; Mahloujifar et al., 2022; Balle et al., 2022; Bhowmick et al., 2018; Stock et al., 2022; Balle et al., 2022; Guo et al., 2022; Kaissis et al., 2023; 2024), using this connection to audit differential privacy (Steinke et al., 2024a; Pillutla et al., 2024; Jagielski et al., 2020; Ding et al., 2018; Bichsel et al., 2018; Nasr et al., 2021; 2023; Steinke et al., 2024b; Tramer et al., 2022; Bichsel et al., 2021; Lu et al., 2022; Andrew et al., 2023; Cebere et al., 2024; Chadha et al., 2024). Some studies have investigated empirical methods to prevent membership

inference attacks that do not rely on differential privacy (Hyland & Tople, 2019; Jia et al., 2019; Chen & Pattabiraman, 2023; Li et al., 2024; Tang et al., 2022; Nasr et al., 2018). An intriguing avenue for future research is to use the concept of empirical privacy to compare the performance of these empirical methods with provable methods, such as DP-SGD.

Algorithm 1 Membership inference in one run game

Input: Oracle access to a mechanism $\mathcal{M}(\cdot)$, A training dataset \mathcal{T} , An indexed canary set $\mathcal{C} = \{x_i; i \in [m]\}$, An attack algorithm A .

- 1: Set $m = |\mathcal{C}|$
 - 2: Sample $u = (u_1, \dots, u_m) \sim \text{Bernoulli}(0.5)^m$, a binary vector where $u_i = 1$ with probability 0.5.
 - 3: Let $\mathcal{S} = \{\mathcal{C}[u_i]; u_i = 1\}_{i \in [m]}$, the subset of selected elements in \mathcal{C} .
 - 4: Run mechanism M on $\mathcal{T} \cup \mathcal{S}$ to get output θ .
 - 5: Run membership inference attack A on θ to get set of membership predictions $v = (v_1, \dots, v_m)$ which is supported on $\{0, 1, \perp\}^m$.
 - 6: Count c , the number of correct guesses where $u_i = v_i$ and c' the total number of guesses where $v_i \neq \perp$.
 - 7: **return** (c, c') .
-

Auditing by reconstruction: We also propose an alternative way to perform auditing by reconstruction attacks. This setup starts with a training set \mathcal{S}_t , similar to the membership inference setting. Then, we have a family of m canary sets $\{\mathcal{S}_c^i; i \in [m]\}$ where each \mathcal{S}_c^i contains k distinct examples. Before training, we construct a set \mathcal{S}_s of size m by uniformly sampling an example from each \mathcal{S}_c^i . Then, the adversary tries to find out which examples were sampled from each canary set \mathcal{S}_c^i by inspecting the model. We recognize that this might be different from what one may consider a true “reconstruction attack”, because the adversary is only performing a selection. However, if you consider the set size to be arbitrary large, and the distribution on the set to be arbitrary, then this will be general enough to cover various notions of reconstruction. We also note that Hayes et al. (2023) use the same setup to measure the performance of the reconstruction attacks.

Algorithm 2 Reconstruction in one run game

Input: Oracle access to a mechanism $\mathcal{M}(\cdot)$, A training dataset \mathcal{T} , number of canaries m , number of options for each canary k , a matrix of canaries $\mathcal{C} = \{x_j^i\}_{i \in [m], j \in [k]}$, an attack algorithm A .

- 1: Let $u = (u_1, \dots, u_m)$ be a vector uniformly sampled from $[k]^m$.
 - 2: Let $\mathcal{S} = \{x_{u_i}^i\}_{i \in [m]}$.
 - 3: Run mechanism \mathcal{M} on $\mathcal{S} \cup \mathcal{T}$ to get output θ .
 - 4: Run a reconstruction attack A on θ to get a vector $v = (v_1, \dots, v_m)$ which is a vector in $([k] \cup \{\perp\})^m$.
 - 5: Count c the number of coordinates where $u_i = v_i$ and c' the number of coordinates where $v_i \neq \perp$.
 - 6: **return** (c, c') .
-

3 IMPLICATIONS OF f -DP FOR GUESSING GAMES

In this section, we explore the implications of f -DP for guessing games. Specifically, we focus on bounding the probability of making more than c correct guesses for adversaries that make at most c' guesses. We begin by stating our main theorem, followed by an explanation of how it can be applied to audit the privacy of a mechanism.

Theorem 9. [Bounds for adversary with bounded guesses] Let $M : [k]^m \rightarrow \Theta$ be a f -DP mechanism. Let \mathbf{u} be a random variable uniformly distributed on $[k]^m$. Let $A : \Theta \rightarrow ([k] \cup \{\perp\})^m$ be a guessing

adversary which always makes at most c' guesses, that is

$$\forall \theta \in \Theta, \Pr \left[\left(\sum_{i=1}^m I(A(\theta)_i \neq \perp) \right) > c' \right] = 0,$$

and let $\mathbf{v} \equiv A(M(\mathbf{u}))$. Define $p_i = \Pr \left[\left(\sum_{j \in [m]} \mathbf{I}(\mathbf{u}_j = \mathbf{v}_j) \right) = i \right]$. For all subset of indices $T \subseteq [c']$, we have

$$\sum_{i \in T} \frac{i}{m} p_i \leq \bar{f} \left(\frac{1}{k-1} \sum_{i \in T} \frac{c' - i + 1}{m} p_{i-1} \right).$$

This Theorem, which we consider to be our main technical contribution, provides a nice invariant that bounds the probability p_i (probability of making exactly i correct guesses) based on the value of other p_j s. Imagine P_f to be a set of vectors $p = (p_1, \dots, p_{c'})$ that could be realized for an attack on a f -DP mechanism. Theorem 9 significantly confines this set. However, this still does not resolve the auditing task. We are interested in bounding $\max_{p \in P_f} \sum_{i=c}^{c'} p_i$, the maximum probability that an adversary can make more than c correct guesses for an f -DP mechanism. Next, we show how we can algorithmically leverage the limitations imposed by Theorem 9 and calculate an upper bound on $\max_{p \in P_f} \sum_{i=c}^{c'} p_i$.

3.1 NUMERICALLY BOUNDING THE TAIL

In this subsection, we specify our procedure for bounding the tail of the distribution and hence the accuracy of our auditing procedure. Our algorithm needs oracle access to f and \bar{f} and decides an upper bound on the probability of an adversary making c correct guesses in a guessing game with alphabet size k and a mechanism that satisfies f -DP. This algorithm relies on the confinement imposed by Theorem 9. Note that Algorithm 3 is a decision algorithm, it takes a value τ and decide if the probability of making more than c correct guesses is less than or equal to τ . We can turn this algorithm to a estimation algorithm by performing a binary search on the value of τ . However, for our use cases, we are interested in a fixed τ . This is because we (similar to (Steinke et al., 2023)) want to set the accuracy of our audit to be a fixed value such as 0.95.

Algorithm 3 Numerically deciding an upper bound probability of making more than c correct guesses

Input: Oracle access to \bar{f} and \bar{f}^{-1} , number of guesses c' , number of correct guesses c , number of samples m , alphabet size k , probability threshold τ (default is $\tau = 0.05$).

```

1:  $\forall 0 \leq i \leq c$  set  $h[i] = 0$ , and  $r[i] = 0$ .
2: set  $r[c] = \tau \cdot \frac{c}{m}$ .
3: set  $h[c] = \tau \cdot \frac{c'-c}{m}$ .
4: for  $i \in [c-1, \dots, 0]$  do
5:    $h[i] = (k-1)\bar{f}^{-1}(r[i+1])$ 
6:    $r[i] = r[i+1] + \frac{i}{c'-i} \cdot (h[i] - h[i+1])$ .
7: end for
8: if  $r[0] + h[0] \geq \frac{c'}{m}$  then
9:   Return True; (Probability of  $c$  correct guesses (out of  $c'$ ) is less than  $\tau$ ).
10: else
11:   Return False; (Probability of having  $c$  correct guesses (out of  $c'$ ) could be more than  $\tau$ ).
12: end if

```

Theorem 10. If Algorithm 3 returns True on inputs \bar{f}, k, m, c, c' and τ , then for any f -DP mechanism $M: [k]^m \rightarrow \Theta$, any guessing adversary $A: \Theta \rightarrow ([k] \cup \{\perp\})^m$ with at most c' guesses, defining \mathbf{u} to be uniform over $[k]^m$, and setting $\mathbf{v} \equiv A(M(\mathbf{u}))$, we have $\Pr[(\sum_{i=1}^m \mathbf{I}(\mathbf{u}_i = \mathbf{v}_i)) \geq c] \leq \tau$.

In a nutshell, this algorithm tries to obtain an upper bound on the sum $p_c + p_{c+1} + \dots, p_{c'}$. We assume this probability is greater than τ , and we obtain lower bound on $p_{c-1} + p_c + \dots + p_{c'}$ based on this assumption. We keep doing this recursively until we have a lower bound on $p_0 + \dots + p_{c'}$. If this lower bound is greater than 1, then we have a contradiction and we return true. The detailed

proof of this Theorem is involved and requires careful analysis. We defer the full proof of Theorem to appendix.

Auditing f -DP with Algorithm 3: When auditing the f -DP for a mechanism, we assume we have injected m canaries, and ran an adversary that is allowed to make c' guesses and recorded that the adversary have made c correct guesses. In such scenario, we will reject the hypothesized privacy of the mechanism if the probability of this observation is less than a threshold τ , which we by default set to 0.05. To this end, we just call Algorithm 3 with parameters $c, c', m, \tau = 0.05$ and f . Then if the algorithm returns *True*, we will reject the privacy hypothesis and approve it otherwise.

Empirical privacy: Although auditing in essence is a hypothesis testing, previous work has used auditing algorithms to calculate empirical privacy as defined in definition 7. In this work, we follow the same route. For simplicity, we only consider an ordered set of privacy hypotheses h_1, \dots, h_w as our family of f -DP curves. These sets are ordered in their strength, meaning that any mechanism that satisfies h_i , would also satisfy h_j for all $j < i$. Then, we would report the strongest privacy hypothesis that passes the test as the empirical privacy of the mechanism.

4 EXPERIMENTS

Most of our experiments are conducted in an *idealized setting*, similar to that used in Steinke et al. (2023), unless otherwise stated. In this setting, the attack success rate is automatically calculated to simulate the expected number of correct guesses by an optimal adversary (Details of the idealized setting are provided in Algorithm 4 in Appendix). We then use this expected number as the default value for the number of correct guesses to derive the empirical ϵ . More specifically, as specified in Definition 6, we instantiate our auditing with a game and evaluation setup. We use Algorithm 4 in Appendix as our game setup. This algorithm returns the number of guesses and the number of correct guesses as the observation from the game. Then, we use Algorithm 3 as our evaluation setup to audit an f -DP curve based on the observation from Algorithm 4. Note that in our comparison with the auditing of Steinke et al., we always use the same membership inference game setup ($k = 2$) as defined in their work. This ensures that our comparison is only on the evaluation part of the audit procedure.

In all experiments, we use empirical ϵ as the primary metric for evaluating our bounds. As described in Section 3.1, we need an ordered set of f -DP curves to obtain empirical privacy. In our experiments, we use f -DP curves for Gaussian mechanisms with varying standard deviations (this forms an ordered set because the f -DP curve of a Gaussian mechanism with a higher standard deviation dominates that of a lower standard deviation). For sub-sampled Gaussian mechanisms, the ordered set consists of f -DP curves for sub-sampled Gaussian mechanisms with the given sub-sampling rate and number of steps and different noise standard deviations.

4.1 COMPARISON WITH STEINKE ET AL. (2023)

In this section, we evaluate our auditing method for membership inference in an idealized setting, using the work of Steinke et al. (2023) as our main baseline. We compare our approach directly to their work, which operates in the same setting as ours.

Simple Gaussian Mechanism: In the first experiment (Figure 1), we audit a simple Gaussian mechanism, varying the standard deviations from $[0.5, 1.0, 2.0, 4.0]$, resulting in different theoretical ϵ values. We vary the number of canaries (m) from 10^2 to 10^7 for auditing, set the bucket size to $k = 2$, and adjust the number of guesses (c') for each number of canaries. For each combination of m, c' , and each standard deviation, we calculate the expected number of correct guesses (c) using Algorithm 4 (the idealized setting). We then audit all tuples of (m, c, c') using the f -DP curves of the Gaussian mechanism, selecting the c that achieves the highest empirical ϵ as the reported empirical ϵ for m canaries at a given standard deviation.

We also apply the same setup for the auditing procedure of Steinke et al. (2023), differing only in the way empirical privacy is calculated. Figure 1 demonstrates that our approach outperforms the empirical privacy results from Steinke et al. Interestingly, while the bound in Steinke et al. (2023) degrades as the number of canaries increases, our bounds continue to improve.

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

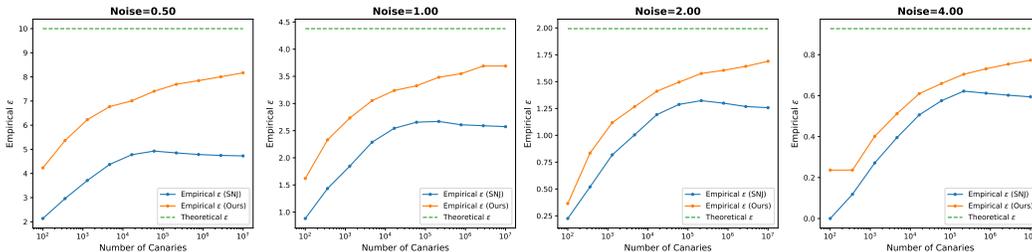


Figure 1: Comparison between our empirical privacy lower bounds and that of Steinke et al. (2023)

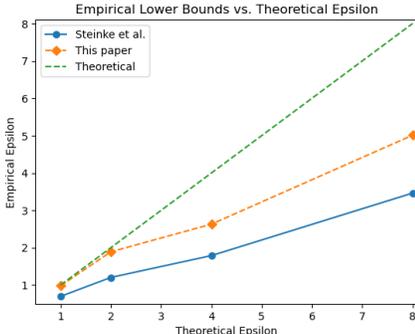


Figure 2: Comparison with auditing procedure of Steinke et al. (2023) on auditing CIFAR-10 in white-box setting using gradient-based membership inference attacks.

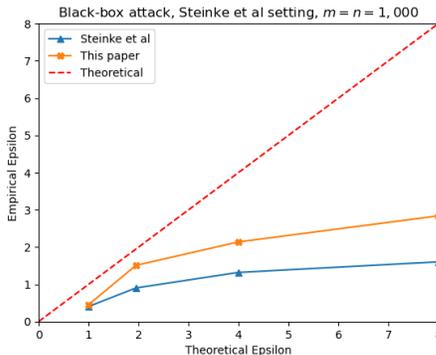


Figure 3: Comparison with auditing procedure of Steinke et al. (2023) on auditing CIFAR-10 in black-box setting.

Experiments on CIFAR-10: We also run experiments on CIFAR-10 on a modified version of the WRN16-4 (Zagoruyko & Komodakis, 2016) architecture, which substitutes batch normalization with group normalization. We follow the setting proposed by Sander et al. (2023), which use custom augmentation multiplicity (i.e., random crop around the center with 20 pixels padding with reflect, random horizontal flip and jitter) and apply an exponential moving average of the model weights with a decay parameter of 0.9999. We run white-box membership inference attacks by following the strongest attack used in the work of Steinke et al. (2023), where the auditor injects multiple canaries in the training set with crafted gradients. More precisely, each canary gradient is set to zero except at a single random index (“Dirac canary” Nasr et al. (2023)). Note that in the white-box attack, the auditor has access to all intermediate iterations of DP-SGD. The attack scores are computed as the dot product between the gradient update during consecutive model iterates and the aggregated gradients from dp-sgd. As done in the work of Steinke et al. (2023), we audit CIFAR-10 model with $m = 5,000$ canaries and all training points from CIFAR-10 $n = 50,000$ for the attack. We set the batch size to 4,096, use augmented multiplicity of $K = 16$ and train for 2,500 DP-SGD steps. For $\epsilon = 8.0, \delta = 10^{-5}$, we achieved 77% accuracy when auditing, compared to 80% without injected canaries. Figure 2 shows the comparison between the auditing scheme by Steinke et al. (2023) with ours for different values of theoretical ϵ . We are able to achieve tighter empirical lower bounds. We also report the performance of the black-box attack, where the auditor does not control the training pipeline and can only compute memberships scores (losses) from the final model. Figure 3 shows how we are able to achieve tighter lower bounds compared to Steinke et al. (2023) where we set $m = 1,000$ and all training samples are used for auditing ($m = n$). This corresponds to the stronger setup for the black-box auditor in Steinke et al. (2023).

Finally, we report the results of auditing the robust membership inference attack Zarifzadeh et al. (2023) (RMIA), which to the best of our knowledge represents the State-of-The-Art (SoTA) black-box membership inference attack on CIFAR-10 from the literature. We reproduce the results in Zarifzadeh et al. (2023) with a non-private WideResNet model (with depth 28 and width 2) for 100 training epochs on half of the dataset chosen at random resulting on a test accuracy of 92.2%. We run the low-cost black-box membership inference attack using 2 reference models

in the offline setting Zarifzadeh et al. (2023). We audit with $m = 5,000$ canaries and report in Figure 4 the comparison between our scheme and Steinke et al. (2023) with different abstention values. Our auditing method clearly outperforms Steinke et al. for all bounded guesses settings, with higher empirical epsilon for larger abstentions values (i.e., smaller number of guesses).

Why is our bound better than Steinke et al. (2023)?

The bounds in Steinke et al. audit approximate DP. That is, they take DP parameters (ϵ, δ) and prove an upper bound on the probability of any adversary obtaining c' correct guesses out of c total guesses, given m canaries available. For the case of $\delta = 0$, their bound is tight. For the case of $\delta > 0$, however, they need to define a set of undesirable events and bound their collective probability. This incurs an additional $O(m \cdot \delta)$ in the probability. The reason why their bounds start to degrade when we increase m is this very fact. The $m \cdot \delta$ term starts to dominate and causes the empirical epsilon estimation to become worse. The reason we do not observe this behavior is that we do not use (ϵ, δ) to approximate the privacy curve, we use the exact curve as is. As we know, the linear approximation of privacy curve is optimal only in a single point for mechanisms that we are interested in (e.g. the Gaussian mechanism). Namely, there is only a single probability $p' \in [0, 1]$ where we have

$$p = \Pr[M(D) \in E] \quad \text{and} \quad e^\epsilon \cdot p + \delta = \Pr[M(D') \in E].$$

Our bound is designed to avoid this issue. We derive a bound that uses the exact f -DP curve, which ensures that for all probabilities $p \in [0, 1]$ the upper bound on the blow-up of events of size p is tight. Moreover, the way we invoke our Theorem 9 in our numerical estimation 3 is designed to apply the bound on events that can be simultaneously tight. This way, our bound does not have the problem of getting worse as the number of samples increases.

Note that this does not mean that there is no way to improve our bound. We still see some gap between the empirical epsilon and the true epsilon. The reason for this, we believe, is in the way numerical tail bound in Algorithm 10 is designed. In this algorithm, we make some relaxations that can be a source of sub-optimality. Specifically, our analysis benefits from the fact that the expectation of correct guesses, conditioned on the correct guesses being greater than c divided by the expectation incorrect guesses conditioned on the same event is greater than c/c' . This step is not tight as we cannot have a mechanism where the adversary makes exactly c correct guesses with probability greater than 0, while making more than c correct guesses with probability exactly 0. For a more interested reader, Equations 6 and 7 in the proof of Theorem 10 is a source of sub-optimality that future work can resolve.

5 CONCLUSIONS AND LIMITATIONS

We introduce a new approach for auditing the privacy of algorithms in a single run using f -DP curves. This method enables more accurate approximations of the true privacy guarantees, addressing the risk of a "false sense of privacy" that may arise from previous approximation techniques. By leveraging the entire f -DP curve, rather than relying solely on point estimates, our approach provides a more nuanced understanding of privacy trade-offs. This allows practitioners to make more informed decisions regarding privacy-utility trade-offs in real-world applications. However, our approach does not provide a strict upper bound on privacy guarantees but instead offers an estimate of the privacy parameters that can be expected in practical scenarios. We also recognize that, despite the improvements over prior work, we still observe a gap between the empirical and theoretical privacy reported in the "one run" setting. Future work could focus on closing this gap to further enhance the reliability of empirical privacy estimations.

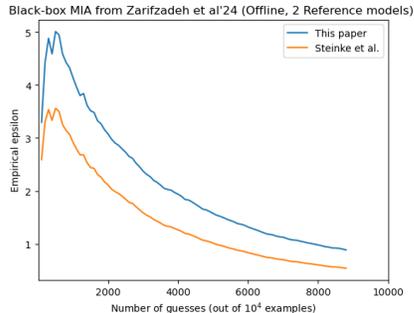


Figure 4: Comparison with auditing procedure of Steinke et al. (2023) on non-private model trained on CIFAR-10 against black-box RMIA method Zarifzadeh et al. (2023).

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

REFERENCES

- 540
541
542 Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and
543 Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC*
544 *conference on computer and communications security*, pp. 308–318, 2016.
- 545 Galen Andrew, Peter Kairouz, Sewoong Oh, Alina Oprea, H Brendan McMahan, and Vinith
546 Suriyakumar. One-shot empirical privacy estimation for federated learning. *arXiv preprint*
547 *arXiv:2302.03098*, 2023.
- 548 Borja Balle, Giovanni Cherubin, and Jamie Hayes. Reconstructing training data with informed
549 adversaries. In *2022 IEEE Symposium on Security and Privacy (SP)*, pp. 1138–1156. IEEE, 2022.
- 550
551 Martin Bertran, Shuai Tang, Aaron Roth, Michael Kearns, Jamie H Morgenstern, and Steven Z Wu.
552 Scalable membership inference attacks via quantile regression. *Advances in Neural Information*
553 *Processing Systems*, 36, 2024.
- 554 Abhishek Bhowmick, John Duchi, Julien Freudiger, Gaurav Kapoor, and Ryan Rogers. Protec-
555 tion against reconstruction and its applications in private federated learning. *arXiv preprint*
556 *arXiv:1812.00984*, 2018.
- 557 Benjamin Bichsel, Timon Gehr, Dana Drachler-Cohen, Petar Tsankov, and Martin Vechev. Dp-finder:
558 Finding differential privacy violations by sampling and optimization. In *Proceedings of the 2018*
559 *ACM SIGSAC Conference on Computer and Communications Security*, pp. 508–524, 2018.
- 560 Benjamin Bichsel, Samuel Steffen, Ilija Bogunovic, and Martin Vechev. Dp-sniper: Black-box
561 discovery of differential privacy violations using classifiers. In *2021 IEEE Symposium on Security*
562 *and Privacy (SP)*, pp. 391–409. IEEE, 2021.
- 563
564 Nicholas Carlini, Steve Chien, Milad Nasr, Shuang Song, Andreas Terzis, and Florian Tramer.
565 Membership inference attacks from first principles. In *2022 IEEE Symposium on Security and*
566 *Privacy (SP)*, pp. 1897–1914. IEEE, 2022.
- 567 Tudor Cebere, Aurélien Bellet, and Nicolas Papernot. Tighter privacy auditing of dp-sgd in the
568 hidden state threat model. *arXiv preprint arXiv:2405.14457*, 2024.
- 569
570 Karan Chadha, Matthew Jagielski, Nicolas Papernot, Christopher Choquette-Choo, and Milad Nasr.
571 Auditing private prediction. *arXiv preprint arXiv:2402.09403*, 2024.
- 572
573 Kamalika Chaudhuri, Claire Monteleoni, and Anand D Sarwate. Differentially private empirical risk
574 minimization. *Journal of Machine Learning Research*, 12(3), 2011.
- 575 Zitao Chen and Karthik Pattabiraman. Overconfidence is a dangerous thing: Mitigating membership
576 inference attacks by enforcing less confident prediction. *arXiv preprint arXiv:2307.01610*, 2023.
- 577
578 Zeyu Ding, Yuxin Wang, Guanhong Wang, Danfeng Zhang, and Daniel Kifer. Detecting violations
579 of differential privacy. In *Proceedings of the 2018 ACM SIGSAC Conference on Computer and*
580 *Communications Security*, pp. 475–489, 2018.
- 581 Jinshuo Dong, Aaron Roth, and Weijie J Su. Gaussian differential privacy. *arXiv preprint*
582 *arXiv:1905.02383*, 2019.
- 583
584 Michael Duan, Anshuman Suri, Niloofar Mireshghallah, Sewon Min, Weijia Shi, Luke Zettlemoyer,
585 Yulia Tsvetkov, Yejin Choi, David Evans, and Hannaneh Hajishirzi. Do membership inference
586 attacks work on large language models? *arXiv preprint arXiv:2402.07841*, 2024.
- 587 Cynthia Dwork. Differential privacy. In *International colloquium on automata, languages, and*
588 *programming*, pp. 1–12. Springer, 2006.
- 589
590 Chuan Guo, Brian Karrer, Kamalika Chaudhuri, and Laurens van der Maaten. Bounding training
591 data reconstruction in private (deep) learning. In *International Conference on Machine Learning*,
592 pp. 8056–8071. PMLR, 2022.
- 593
594 Jamie Hayes, Saeed Mahloujifar, and Borja Balle. Bounding training data reconstruction in dp-sgd.
arXiv preprint arXiv:2302.07225, 2023.

- 594 Nils Homer, Szabolcs Szelinger, Margot Redman, David Duggan, Waibhav Tembe, Jill Muehling,
595 John V Pearson, Dietrich A Stephan, Stanley F Nelson, and David W Craig. Resolving individuals
596 contributing trace amounts of dna to highly complex mixtures using high-density snp genotyping
597 microarrays. *PLoS genetics*, 4(8):e1000167, 2008.
- 598 Hongsheng Hu, Zoran Salcic, Lichao Sun, Gillian Dobbie, Philip S Yu, and Xuyun Zhang. Member-
599 ship inference attacks on machine learning: A survey. *ACM Computing Surveys (CSUR)*, 54(11s):
600 1–37, 2022.
- 601 Stephanie L Hyland and Shruti Tople. On the intrinsic privacy of stochastic gradient descent. *Preprint*
602 *at <https://arxiv.org/pdf/1912.02919.pdf>*, 2019.
- 603 Matthew Jagielski, Jonathan Ullman, and Alina Oprea. Auditing differentially private machine
604 learning: How private is private sgd? *Advances in Neural Information Processing Systems*, 33:
605 22205–22216, 2020.
- 606 Jinyuan Jia, Ahmed Salem, Michael Backes, Yang Zhang, and Neil Zhenqiang Gong. Memguard: De-
607 fending against black-box membership inference attacks via adversarial examples. In *Proceedings*
608 *of the 2019 ACM SIGSAC conference on computer and communications security*, pp. 259–274,
609 2019.
- 610 Georgios Kaissis, Jamie Hayes, Alexander Ziller, and Daniel Rueckert. Bounding data recon-
611 struction attacks with the hypothesis testing interpretation of differential privacy. *arXiv preprint*
612 *arXiv:2307.03928*, 2023.
- 613 Georgios Kaissis, Alexander Ziller, Stefan Kolek, Anneliese Riess, and Daniel Rueckert. Optimal
614 privacy guarantees for a relaxed threat model: Addressing sub-optimal adversaries in differentially
615 private machine learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- 616 Klas Leino and Matt Fredrikson. Stolen memories: Leveraging model memorization for calibrated
617 {White-Box} membership inference. In *29th USENIX security symposium (USENIX Security 20)*,
618 pp. 1605–1622, 2020.
- 619 Jiacheng Li, Ninghui Li, and Bruno Ribeiro. {MIST}: Defending against membership inference
620 attacks through {Membership-Invariant} subspace training. In *33rd USENIX Security Symposium*
621 *(USENIX Security 24)*, pp. 2387–2404, 2024.
- 622 Fred Lu, Joseph Munoz, Maya Fuchs, Tyler LeBlond, Elliott Zaresky-Williams, Edward Raff, Francis
623 Ferraro, and Brian Testa. A general framework for auditing differentially private machine learning.
624 *Advances in Neural Information Processing Systems*, 35:4165–4176, 2022.
- 625 Saeed Mahloujifar, Alexandre Sablayrolles, Graham Cormode, and Somesh Jha. Optimal membership
626 inference bounds for adaptive composition of sampled gaussian mechanisms. *arXiv preprint*
627 *arXiv:2204.06106*, 2022.
- 628 Jagielski Matthew, Nasr Milad, Choquette-Choo Christopher, Lee Katherine, and Carlini Nicholas.
629 Students parrot their teachers: Membership inference on model distillation. *arXiv preprint arXiv:*
630 *2303.03446*, 2023.
- 631 Ilya Mironov. Rényi differential privacy. In *2017 IEEE 30th computer security foundations symposium*
632 *(CSF)*, pp. 263–275. IEEE, 2017.
- 633 Milad Nasr, Reza Shokri, and Amir Houmansadr. Machine learning with membership privacy using
634 adversarial regularization. In *Proceedings of the 2018 ACM SIGSAC conference on computer and*
635 *communications security*, pp. 634–646, 2018.
- 636 Milad Nasr, Shuang Songi, Abhradeep Thakurta, Nicolas Papernot, and Nicholas Carlin. Adversary
637 instantiation: Lower bounds for differentially private machine learning. In *2021 IEEE Symposium*
638 *on security and privacy (SP)*, pp. 866–882. IEEE, 2021.
- 639 Milad Nasr, Jamie Hayes, Thomas Steinke, Borja Balle, Florian Tramèr, Matthew Jagielski, Nicholas
640 Carlini, and Andreas Terzis. Tight auditing of differentially private machine learning. *arXiv*
641 *preprint arXiv:2302.07956*, 2023.

- 648 Krishna Pillutla, Galen Andrew, Peter Kairouz, H Brendan McMahan, Alina Oprea, and Sewoong Oh.
649 Unleashing the power of randomization in auditing differentially private ml. *Advances in Neural*
650 *Information Processing Systems*, 2024.
- 651 Alexandre Sablayrolles, Matthijs Douze, Cordelia Schmid, Yann Ollivier, and Hervé Jégou. White-
652 box vs black-box: Bayes optimal strategies for membership inference. In *International Conference*
653 *on Machine Learning*, pp. 5558–5567. PMLR, 2019.
- 654 Tom Sander, Pierre Stock, and Alexandre Sablayrolles. Tan without a burn: Scaling laws of dp-sgd.
655 In *International Conference on Machine Learning*. PMLR, 2023.
- 656 Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks
657 against machine learning models. In *2017 IEEE symposium on security and privacy (SP)*, pp. 3–18.
658 IEEE, 2017.
- 659 Thomas Steinke, Milad Nasr, and Matthew Jagielski. Privacy auditing with one (1) training run.
660 *arXiv preprint arXiv:2305.08846*, 2023.
- 661 Thomas Steinke, Milad Nasr, Arun Ganesh, Borja Balle, Christopher A Choquette-Choo, Matthew
662 Jagielski, Jamie Hayes, Abhradeep Guha Thakurta, Adam Smith, and Andreas Terzis. The last
663 iterate advantage: Empirical auditing and principled heuristic analysis of differentially private sgd.
664 *arXiv preprint arXiv:2410.06186*, 2024a.
- 665 Thomas Steinke, Milad Nasr, and Matthew Jagielski. Privacy auditing with one (1) training run.
666 *Advances in Neural Information Processing Systems*, 36, 2024b.
- 667 Pierre Stock, Igor Shilov, Ilya Mironov, and Alexandre Sablayrolles. Defending against reconstruction
668 attacks with r -enyi differential privacy. *arXiv preprint arXiv:2202.07623*, 2022.
- 669 Xinyu Tang, Saeed Mahloujifar, Liwei Song, Virat Shejwalkar, Milad Nasr, Amir Houmansadr, and
670 Prateek Mittal. Mitigating membership inference attacks by {Self-Distillation} through a novel
671 ensemble architecture. In *31st USENIX Security Symposium (USENIX Security 22)*, pp. 1433–1450,
672 2022.
- 673 Florian Tramer, Andreas Terzis, Thomas Steinke, Shuang Song, Matthew Jagielski, and Nicholas
674 Carlini. Debugging differential privacy: A case study for privacy auditing. *arXiv preprint*
675 *arXiv:2202.12219*, 2022.
- 676 Jiachen T Wang, Saeed Mahloujifar, Tong Wu, Ruoxi Jia, and Prateek Mittal. A randomized approach
677 for tight privacy accounting. *arXiv preprint arXiv:2304.07927*, 2023.
- 678 Lauren Watson, Chuan Guo, Graham Cormode, and Alex Sablayrolles. On the importance of difficulty
679 calibration in membership inference attacks. *arXiv preprint arXiv:2111.08440*, 2021.
- 680 Jiayuan Ye and Reza Shokri. Differentially private learning needs hidden state (or much faster
681 convergence). *Advances in Neural Information Processing Systems*, 35:703–715, 2022.
- 682 Jiayuan Ye, Aadyaa Maddi, Sasi Kumar Murakonda, Vincent Bindschaedler, and Reza Shokri.
683 Enhanced membership inference attacks against machine learning models. In *Proceedings of the*
684 *2022 ACM SIGSAC Conference on Computer and Communications Security*, pp. 3093–3106, 2022.
- 685 Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*,
686 2016.
- 687 Sajjad Zarifzadeh, Philippe Cheng-Jie Marc Liu, and Reza Shokri. Low-cost high-power membership
688 inference by boosting relativity. 2023.
- 689 Yuqing Zhu, Jinshuo Dong, and Yu-Xiang Wang. Optimal accounting of differential privacy via
690 characteristic function. In *International Conference on Artificial Intelligence and Statistics*, pp.
691 4782–4817. PMLR, 2022.
- 692
693
694
695
696
697
698
699
700
701

702 A PROOFS

703 A.1 PROOF OUTLINE FOR THEOREM 9

704 In this subsection, we outline the main ingredients we need to prove our Theorem 9. We also provide
705 the full proof for a simplified version of Theorem 9 using these ingredients. First, we have a Lemma
706 that bounds the probability of any event conditioned on correctly guessing a single canary.

707 **Lemma 11.** *Let $M : [k]^m \rightarrow \Theta$ be a mechanism that satisfies f -DP. Also let $A : \Theta \rightarrow ([k] \cup \{\perp\})^m$
708 be a guessing attack. Let \mathbf{u} be a random variable uniformly distributed over $[k]^m$ and let $\mathbf{v} \equiv$
709 $A(M(\mathbf{u}))$. Then for any subset $E \subseteq \Theta$ we have*

$$710 f_k'' \left(\Pr [M(\mathbf{u}) \in E] \right) \leq \Pr [M(\mathbf{u}) \in E \text{ and } u_1 = v_1] \leq f_k' \left(\Pr [M(\mathbf{u}) \in E] \right)$$

711 where

$$712 f_k'(x) = \sup\{\alpha; \alpha + f\left(\frac{x - \alpha}{k - 1}\right) \leq 1\} \text{ and } f_k''(x) = \inf\{\alpha; (k - 1)f(\alpha) + x - \alpha \leq 1\}.$$

713 This Lemma which is a generalization and an improvement over the main Theorem of (Hayes et al.,
714 2023), shows that the probability of an event cannot change too much if we condition on the success
715 of adversary on one of the canaries. Note that this Lemma immediately implies a bound on the
716 expected number of correct guesses by any guessing adversary (by just using linearity of expectation).
717 However, here we are not interested in expectations. Rather, we need to derive tail bounds. The proof
718 of Theorem 9 relies on some key properties of the f' and f'' functions defined in the statement of
719 Lemma 11. These properties are specified in the following Proposition and proved in the Appendix.

720 **Proposition 12.** *The functions f_k' as defined in Lemma 11 is increasing and concave. The function
721 f_k'' as defined in Lemma 11 is increasing and convex.*

722 Now, we are ready to outline the proof of a simplified variant of our Theorem 9 for adversaries that
723 make a guess on all canaries. This makes the proof much simpler and enables us to focus more on the
724 key steps in the proof.

725 **Theorem 13** (Special case of 9). *Let $M : [k]^m \rightarrow \Theta$ be a f -DP mechanism. Let \mathbf{u} be a random
726 variable uniformly distributed on $[k]^m$. Let $A : \Theta \rightarrow [k]^m$ be a guessing adversary and let $\mathbf{v} \equiv$
727 $A(M(\mathbf{u}))$. Define $p_i = \Pr \left[\left(\sum_{j \in [m]} \mathbf{I}(\mathbf{u}_j = \mathbf{v}_j) \right) = i \right]$. For all subset of indices $T \subseteq [m]$, we have*

$$728 \sum_{i \in T} \frac{i}{m} p_i \leq \bar{f} \left(\frac{1}{k - 1} \sum_{i \in T} \frac{m - i + 1}{m} p_{i-1} \right)$$

729 *Proof.* Let us define a random variable $\mathbf{t} = (\mathbf{t}_1, \dots, \mathbf{t}_m)$ which is defined as $\mathbf{t}_i = \mathbf{I}(\mathbf{u}_i = \mathbf{v}_i)$ We
730 have

$$731 p_c = \Pr \left[\sum_{i=1}^m \mathbf{t}_i = c \right] = \Pr \left[\sum_{i=2}^m \mathbf{t}_i = c - 1 \text{ and } \mathbf{t}_1 = 1 \right] + \Pr \left[\sum_{i=2}^m \mathbf{t}_i = c \text{ and } \mathbf{t}_1 = 0 \right]$$

732 Now by Lemma 11 we have $\Pr \left[\sum_{i=2}^m \mathbf{t}_i = c - 1 \text{ and } \mathbf{t}_1 = 1 \right] \leq f_k' \left(\sum_{i=2}^m \mathbf{t}_i = c - 1 \right)$. This is a nice
733 invariant that we can use but $\sum_{i=2}^m \mathbf{t}_i = c - 1$ could be really small depending on how large m is. To
734 strengthen the bound we sum all p_c 's for $c \in T$, and then apply the lemma on the aggregate. That is

$$735 \sum_{j \in T} p_j = \sum_{j \in T} \Pr \left[\sum_{i=1}^m \mathbf{t}_i = j \right] = \sum_{j \in T} \Pr \left[\sum_{i=2}^m \mathbf{t}_i = j \text{ and } \mathbf{t}_1 = 0 \right] + \sum_{j \in T} \Pr \left[\sum_{i=2}^m \mathbf{t}_i = j - 1 \text{ and } \mathbf{t}_1 = 1 \right]$$

$$736 = \Pr \left[\sum_{i=2}^m \mathbf{t}_i \in T \text{ and } \mathbf{t}_1 = 0 \right] + \Pr \left[1 + \sum_{i=2}^m \mathbf{t}_i \in T \text{ and } \mathbf{t}_1 = 1 \right]$$

Now we only use the inequality from Lemma 11 for the second quantity above. Using the inequality for both probabilities is not ideal because they cannot be tight at the same time. So we have,

$$\sum_{j \in T} p_j \leq \Pr[\sum_{i=2}^m \mathbf{t}_i \in T \text{ and } \mathbf{t}_1 = 0] + f'_k(\Pr[1 + \sum_{i=2}^m \mathbf{t}_i \in T]).$$

Now we use a trick to make this cleaner. We use the fact that this inequality is invariant to the order of indices. So we can permute \mathbf{t}_i 's and the inequality still holds. We have,

$$\begin{aligned} \sum_{j \in T} p_j &\leq \mathbb{E}_{\pi \sim \Pi[m]} [\Pr[\sum_{i=2}^m \mathbf{t}_{\pi(i)} \in T \text{ and } \mathbf{t}_{\pi(1)} = 0]] + \mathbb{E}_{\pi \sim \Pi[m]} [f'_k(\Pr[1 + \sum_{i=2}^m \mathbf{t}_{\pi(i)} \in T])] \\ &\leq \mathbb{E}_{\pi \sim \Pi[m]} [\Pr[\sum_{i=2}^m \mathbf{t}_{\pi(i)} \in T \text{ and } \mathbf{t}_{\pi(1)} = 0]] + f'_k(\mathbb{E}_{\pi \sim \Pi[m]} [\Pr[1 + \sum_{i=2}^m \mathbf{t}_{\pi(i)} \in T]]). \end{aligned}$$

Now we perform a double counting argument. Note that when we permute the order $\sum_{i=2}^m \mathbf{t}_{\pi(i)} = j$ and $\mathbf{t}_{\pi(1)} = 0$ counts each instance t_1, \dots, t_m with exactly j non-zero locations, for exactly $(m-j) \times (m-1)!$ times. Therefore, we have

$$\mathbb{E}_{\pi \sim \Pi[m]} [\Pr[\sum_{i=2}^m \mathbf{t}_{\pi(i)} \in T \text{ and } \mathbf{t}_{\pi(1)} = 0]] = \sum_{j \in T} \frac{m-j}{m} p_j.$$

With a similar argument we have,

$$\mathbb{E}_{\pi \sim \Pi[m]} [\Pr[1 + \sum_{i=2}^m \mathbf{t}_{\pi(i)} \in T]] = \sum_{j \in T} \frac{m-j+1}{m} p_{j-1} + \frac{j}{m} p_j.$$

Then, we have

$$\sum_{j \in T} p_j \leq \sum_{j \in T} \frac{m-j}{m} p_j + f'_k(\sum_{j \in T} \frac{j}{m} p_j + \frac{m-j+1}{m} p_{j-1}).$$

And this implies

$$\sum_{j \in T} \frac{j}{m} p_j \leq f'_k(\sum_{j \in T} \frac{j}{m} p_j + \frac{m-j+1}{m} p_{j-1}).$$

And this, by definition of f'_k implies

$$\sum_{j \in T} \frac{j}{m} p_j \leq \bar{f}(\frac{1}{k-1} \sum_{j \in T} \frac{m-j+1}{m} p_{j-1}).$$

□

A.2 FULL PROOFS

Proof of Lemma 11. Let $p = \Pr[M(\mathbf{u}) \in E \text{ and } u_1 = v_1]$ and $q = \Pr[M(\mathbf{u}) \in E]$. We have

$$\begin{aligned}
810 & p = \sum_{i \in [k]} \Pr[M(\mathbf{u}) \in E \text{ and } u_1 = v_1 = i] \\
811 & = \frac{1}{k} \sum_{i \in [k]} \Pr[M(\mathbf{u}) \in E \text{ and } v_1 = i \mid u_1 = i] \\
812 & = \frac{1}{k} \sum_{i \in [k]} \frac{1}{k-1} \left(\sum_{j \in [k] \setminus \{i\}} \Pr[M(\mathbf{u}) \in E \text{ and } v_1 = i \mid u_1 = j] \right) \\
813 & \stackrel{\text{(By definition of } f\text{-DP)}}{\leq} \frac{1}{k} \sum_{i \in [k]} \frac{1}{k-1} \left(\sum_{j \in [k] \setminus \{i\}} 1 - f(\Pr[M(\mathbf{u}) \in E \text{ and } v_1 = i \mid u_1 = j]) \right) \\
814 & \stackrel{\text{(By convexity of } f\text{)}}{\leq} 1 - f \left(\frac{1}{k} \sum_{i \in [k]} \frac{1}{k-1} \left(\sum_{j \in [k] \setminus \{i\}} \Pr[M(\mathbf{u}) \in E \text{ and } v_1 = i \mid u_1 = j] \right) \right) \\
815 & = 1 - f \left(\frac{1}{k-1} \sum_{i \in [k]} \left(\sum_{j \in [k] \setminus \{i\}} \frac{1}{k} \Pr[M(\mathbf{u}) \in E \text{ and } v_1 = i \mid u_1 = j] \right) \right) \\
816 & = 1 - f \left(\frac{1}{k-1} \sum_{i \in [k]} \left(\sum_{j \in [k] \setminus \{i\}} \Pr[M(\mathbf{u}) \in E \text{ and } v_1 = i \text{ and } u_1 = j] \right) \right) \\
817 & = 1 - f \left(\frac{1}{k-1} \Pr[M(\mathbf{u}) \in E \text{ and } u_1 \neq v_1] \right) \\
818 & = 1 - f \left(\frac{q-p}{k-1} \right).
\end{aligned}$$

819 Similarly we have,

$$\begin{aligned}
820 & p = \sum_{i \in [k]} \Pr[M(\mathbf{u}) \in E \text{ and } u_1 = v_1 = i] \\
821 & = \frac{1}{k} \sum_{i \in [k]} \Pr[M(\mathbf{u}) \in E \text{ and } v_1 = i \mid u_1 = i] \\
822 & = \frac{1}{k} \sum_{i \in [k]} \frac{1}{k-1} \left(\sum_{j \in [k] \setminus \{i\}} \Pr[M(\mathbf{u}) \in E \text{ and } v_1 = i \mid u_1 = j] \right) \\
823 & \stackrel{\text{(By definition of } f\text{-DP)}}{\geq} \frac{1}{k} \sum_{i \in [k]} \frac{1}{k-1} \left(\sum_{j \in [k] \setminus \{i\}} f^{-1}(1 - \Pr[M(\mathbf{u}) \in E \text{ and } v_1 = i \mid u_1 = j]) \right) \\
824 & \stackrel{\text{(By convexity of } f\text{)}}{\geq} f^{-1} \left(\frac{1}{k} \sum_{i \in [k]} \frac{1}{k-1} \left(\sum_{j \in [k] \setminus \{i\}} 1 - \Pr[M(\mathbf{u}) \in E \text{ and } v_1 = i \mid u_1 = j] \right) \right) \\
825 & = f^{-1} \left(\frac{1}{k-1} \sum_{i \in [k]} \left(\sum_{j \in [k] \setminus \{i\}} \frac{1}{k} (1 - \Pr[M(\mathbf{u}) \in E \text{ and } v_1 = i \mid u_1 = j]) \right) \right) \\
826 & = f^{-1} \left(\frac{1}{k-1} \sum_{i \in [k]} \left(\sum_{j \in [k] \setminus \{i\}} \Pr[M(\mathbf{u}) \in E \text{ and } v_1 = i \text{ and } u_1 = j] \right) \right) \\
827 & = f^{-1} \left(\frac{1}{k-1} (1 - \Pr[M(\mathbf{u}) \in E \text{ and } u_1 \neq v_1]) \right) \\
828 & = f^{-1} \left(\frac{1-q+p}{k-1} \right).
\end{aligned}$$

829 This implies that,

$$f(p) \cdot (k-1) + q - p \leq 1$$

□

Proof of Proposition 12. The function is increasing simply because f is decreasing. We now prove concavity. Let $\alpha_1 = f_k(x_1)$ and $\alpha_2 = f_k(x_2)$. By definition of f_k we have

$$\alpha_1 + f\left(\frac{x_1 - \alpha_1}{k-1}\right) \leq 1$$

and

$$\alpha_2 + f\left(\frac{x_2 - \alpha_2}{k-1}\right) \leq 1.$$

Averaging these two we get,

$$\frac{\alpha_1 + \alpha_2}{2} + \frac{f\left(\frac{x_1 - \alpha_1}{k-1}\right) + f\left(\frac{x_2 - \alpha_2}{k-1}\right)}{2} \leq 1$$

By convexity of f we have

$$\frac{\alpha_1 + \alpha_2}{2} + f\left(\frac{\frac{x_1 + x_2}{2} - \frac{\alpha_1 + \alpha_2}{2}}{k-1}\right) \leq 1$$

Therefore, by definition of f'_k , we have $f'_k\left(\frac{x_1 + x_2}{2}\right) \geq \frac{\alpha_1 + \alpha_2}{2}$. Similarly, f''_k is increasing just because f is decreasing. And assuming $\alpha_1 = f_k(x_1)$ and $\alpha_2 = f_k(x_2)$ we have

$$f''_k\left(\frac{x_1 + x_2}{2}\right) \leq \frac{\alpha_1 + \alpha_2}{2}$$

which implies f''_k is convex. □

Proof of Theorem 9. Instead of working with an adversary with c' guesses, we assume we have an adversary that makes a guess on all m inputs, however, it also submits a vector $\mathbf{q} \in \{0, 1\}^m$, with exactly c' 1s and $m - c'$ 0s. So the output of this adversary is a vector $\mathbf{v} \in [k]^m$ and a vector $\mathbf{q} \in \{0, 1\}^m$. Then, only correct guesses that are in locations that \mathbf{q} is non-zero is counted. That is, if we define a random variable $\mathbf{t} = (t_1, \dots, t_m)$ as $t_i = \mathbf{I}(u_i = v_i)$ then we have

$$\begin{aligned} p_c &= \Pr\left[\sum_{i=1}^m t_i \cdot \mathbf{q}_i = c\right] \\ &= \Pr\left[\sum_{i=2}^m t_i = c - 1 \text{ and } t_1 = 1 \text{ and } \mathbf{q}_1 = 1\right] + \Pr\left[\sum_{i=2}^m t_i = c \text{ and } t_1 \cdot \mathbf{q}_1 = 0\right] \end{aligned}$$

Now by Lemma 11 we have

$$\Pr\left[\sum_{i=2}^m t_i = c - 1 \text{ and } t_1 = 1 \text{ and } \mathbf{q}_1 = 1\right] \leq f'_k\left(\sum_{i=2}^m t_i = c - 1 \text{ and } \mathbf{q}_1 = 1\right).$$

This is a nice invariant that we can use but $\sum_{i=2}^m t_i = c - 1$ could be really small depending on how large m is. To strengthen the bound we sum all p_c 's for $c \in T$, and then apply the lemma on the aggregate. That is

$$\begin{aligned} \sum_{j \in T} p_j &= \sum_{j \in T} \Pr\left[\sum_{i=1}^m t_i = j\right] \\ &= \sum_{j \in T} \Pr\left[\sum_{i=2}^m t_i = j \text{ and } t_1 \cdot \mathbf{q}_1 = 0\right] + \sum_{j \in T} \Pr\left[\sum_{i=2}^m t_i = j - 1 \text{ and } t_1 = 1 \text{ and } \mathbf{q}_1 = 1\right] \\ &= \Pr\left[\sum_{i=2}^m t_i \in T \text{ and } t_1 \cdot \mathbf{q}_1 = 0\right] + \Pr\left[1 + \sum_{i=2}^m t_i \in T \text{ and } t_1 = 1 \text{ and } \mathbf{q}_1 = 1\right] \end{aligned}$$

Now we only use the inequality from Lemma 11 for the second quantity above. Using the inequality for both probabilities is not ideal because they cannot be tight at the same time. So we have,

$$\sum_{j \in T} p_j \leq \Pr\left[\sum_{i=2}^m \mathbf{t}_i \in T \text{ and } \mathbf{t}_1 \cdot \mathbf{q}_1 = 0\right] + f'_k\left(\Pr\left[1 + \sum_{i=2}^m \mathbf{t}_i \in T \text{ and } \mathbf{q}_1 = 1\right]\right).$$

Now we use a trick to make this cleaner. We use the fact that this inequality is invariant to the order of indices. So we can permute \mathbf{t}_i 's and the inequality still holds. We have,

$$\begin{aligned} \sum_{j \in T} p_j &\leq \mathbb{E}_{\pi \sim \Pi[m]} \left[\Pr\left[\sum_{i=2}^m \mathbf{t}_{\pi(i)} \in T \text{ and } \mathbf{t}_{\pi(1)} \cdot \mathbf{q}_{\pi(1)} = 0\right] \right] + \mathbb{E}_{\pi \sim \Pi[m]} \left[f'_k\left(\Pr\left[1 + \sum_{i=2}^m \mathbf{t}_{\pi(i)} \in T\right]\right) \right] \\ &\leq \mathbb{E}_{\pi \sim \Pi[m]} \left[\Pr\left[\sum_{i=2}^m \mathbf{t}_{\pi(i)} \in T \text{ and } \mathbf{t}_{\pi(1)} = 0\right] \right] + f'_k\left(\mathbb{E}_{\pi \sim \Pi[m]} \left[\Pr\left[1 + \sum_{i=2}^m \mathbf{t}_{\pi(i)} \in T \text{ and } \mathbf{q}_{\pi(1)} = 1\right] \right]\right). \end{aligned}$$

Now we perform a double counting argument. Note that when we permute the order $\sum_{i=2}^m \mathbf{t}_{\pi(i)} = j$ and $\mathbf{t}_{\pi(1)} = 0$ counts each instance t_1, \dots, t_m with exactly j non-zero locations, for exactly $(m-j) \times (m-1)!$ times. Therefore, we have

$$\mathbb{E}_{\pi \sim \Pi[m]} \left[\Pr\left[\sum_{i=2}^m \mathbf{t}_{\pi(i)} \cdot \mathbf{q}_{\pi(i)} \in T \text{ and } \mathbf{t}_{\pi(1)} \cdot \mathbf{q}_{\pi(1)} = 0\right] \right] = \sum_{j \in T} \frac{m-j}{m} p_j.$$

With a similar argument we have,

$$\mathbb{E}_{\pi \sim \Pi[m]} \left[\Pr\left[1 + \sum_{i=2}^m \mathbf{t}_{\pi(i)} \cdot \mathbf{q}_{\pi(i)} \in T \text{ and } \mathbf{q}_{\pi(1)} = 1\right] \right] = \sum_{j \in T} \frac{c' - j + 1}{m} p_{j-1} + \frac{j}{m} p_j.$$

Then, we have

$$\begin{aligned} \sum_{j \in T} p_j &\leq \sum_{j \in T} \frac{m-j}{m} p_j + f'_k\left(\sum_{j \in T} \frac{j}{m} p_j + \frac{c' - j + 1}{m} p_{j-1}\right) \\ &= \sum_{j \in T} \frac{m-j}{m} p_j + f'_k\left(\sum_{j \in T} \frac{j}{m} p_j + \frac{c' - j + 1}{m} p_{j-1}\right). \end{aligned}$$

And this implies

$$\sum_{j \in T} \frac{j}{m} p_j \leq f'_k\left(\sum_{j \in T} \frac{j}{m} p_j + \frac{c' - j + 1}{m} p_{j-1}\right).$$

And this, by definition of f'_k implies

$$\sum_{j \in T} \frac{j}{m} p_j \leq \bar{f}\left(\frac{1}{k-1} \sum_{j \in T} \frac{c' - j + 1}{m} p_{j-1}\right).$$

□

Proof of Theorem 10. To prove Theorem 10, we first state and prove a lemma which is consequence of Theorem 9.

Lemma 14. For all $c \leq c' \in [m]$ let us define

$$\alpha_c = \sum_{i=c}^{c'} \frac{i}{m} p_i \quad \text{and} \quad \beta_c = \sum_{i=c}^{c'} \frac{c' - i}{m} p_i$$

We also define a family of functions $r = \{r_{i,j} : [0, 1] \times [0, 1] \rightarrow [0, 1]\}_{i \leq j \in [m]}$ and $h = \{h_{i,j} : [0, 1] \rightarrow [0, 1]\}$ that are defined recursively as follows.

972 $\forall i \in [m] : r_{i,i}(\alpha, \beta) = \alpha$ and $h_{i,i}(\alpha, \beta) = \beta$ and for all $i < j$ we have

$$973 \quad h_{i,j}(\alpha, \beta) = (k-1)\bar{f}^{-1}\left(r_{i+1,j}(\alpha, \beta)\right)$$

$$974 \quad r_{i,j}(\alpha, \beta) = r_{i+1,j}(\alpha, \beta) + \frac{i}{c'-i}(h_{i,j}(\alpha, \beta) - h_{i+1,j}(\alpha, \beta))$$

975
976
977 Then for all $i \leq j$ we have

$$978 \quad \alpha_i \geq r_{i,j}(\alpha_j, \beta_j) \quad \text{and} \quad \beta_i \geq h_{i,j}(\alpha_j, \beta_j)$$

979
980 Moreover, for $i < j$, $r_{i,j}$ and $h_{i,j}$ are increasing with respect to their first argument and decreasing
981 with respect to their second argument.

982
983 *Proof of Lemma 14.* We prove this by induction on $j - i$. For $j - i = 0$, the statement is trivially
984 correct. We have

$$985 \quad h_{i,j}(\alpha_j, \beta_j) = (k-1)\bar{f}^{-1}(r_{i+1,j}(\alpha_j, \beta_j)).$$

986 By induction hypothesis, we have $r_{i+1,j}(\alpha_j, \beta_j) \leq \alpha_{i+1}$. Therefore we have

$$987 \quad h_{i,j}(\alpha_j, \beta_j) \leq (k-1)\bar{f}^{-1}(\alpha_{i+1}). \quad (1)$$

988 Now by invoking Theorem 9, we have

$$989 \quad \alpha_{i+1} \leq \bar{f}\left(\frac{\beta_i}{k-1}\right).$$

990 Now since \bar{f} is increasing, this implies

$$991 \quad (k-1)\bar{f}^{-1}(\alpha_{i+1}) \leq \beta_i \quad (2)$$

992 Now putting, inequalities 1 and 2 together we have $h_{i,j}(\alpha_j, \beta_j) \leq \beta_i$. This proves the first part of
993 the induction hypothesis for the function h . Also note that $h_{i,j}$ is increasing in its first component
994 and decreasing in the second component by invoking induction hypothesis and the fact that \bar{f}^{-1} is
995 increasing. Now we focus on function $r_{i,j}$. First note that there is an alternative form for $r_{i,j}$ by
996 opening up the recursive relation. Let $\gamma_z = \frac{z}{c'-z} - \frac{z-1}{c'-z+1}$. We have ,

$$997 \quad r_{i,j}(\alpha, \beta) = r_{j,j}(\alpha, \beta) + \frac{i}{c'-i}h_{i,j}(\alpha, \beta) - \frac{j-1}{c'-j+1}h_{j,j}(\alpha, \beta) + \sum_{z=i+1}^{j-1} \gamma_z h_{z,j}(\alpha, \beta)$$

$$998 \quad = r_{j,j}(\alpha, \beta) + \frac{i}{c'-i}h_{i,j}(\alpha, \beta) - \frac{j}{c'-j}h_{j,j}(\alpha, \beta) + \sum_{z=i+1}^j \gamma_z h_{z,j}(\alpha, \beta)$$

$$999 \quad = \alpha - \frac{j}{c'-j}\beta + \frac{i}{c'-i}h_{i,j}(\alpha, \beta) + \sum_{z=i+1}^j \gamma_z h_{z,j}(\alpha, \beta). \quad (3)$$

1000 Now we show that for all i we have

$$1001 \quad \alpha_i = \frac{i}{c'-i}\beta_i + \sum_{z=i+1}^m \gamma_z \beta_z. \quad (4)$$

1002 This is because we have

$$1003 \quad \alpha_i - \frac{i}{c'-i}\beta_i = \sum_{z=i+1}^{c'} \left(\frac{z}{m} - \frac{i(c'-z)}{(c'-i)m}\right)p_z.$$

1004 On the other hand we have

$$1005 \quad \sum_{z=i+1}^m \gamma_z \beta_z = \sum_{z=i+1}^m \left(\sum_{z'=i+1}^z \gamma_{z'}\right) \frac{c'-z}{m} p_z$$

$$1006 \quad = \sum_{z=i+1}^m \left(\frac{z}{c'-z} - \frac{i}{c'-i}\right) \frac{c'-z}{m} p_z$$

$$1007 \quad = \sum_{z=i+1}^m \left(\frac{z}{m} - \frac{i(c'-z)}{(c'-i)m}\right) p_z$$

and this shows that Equation 4 is correct. Therefore for all $i < j$ we have

$$\alpha_i - \alpha_j = \frac{i}{c' - i} \beta_i - \frac{j}{c' - j} \beta_j + \sum_{z=i+1}^j \gamma_z \beta_z$$

Now, using the induction hypothesis for h we have,

$$\alpha_i \geq \alpha_j + \frac{i}{c' - i} h_{i,j}(\alpha_j, \beta_j) - \frac{j}{c' - j} \beta_j + \sum_{z=i+1}^j \gamma_z h_{z,j}(\alpha_j, \beta_j). \quad (5)$$

Now verify that the right hand side of Equation 5 is equal to $r_{i,j}(\alpha_j, \beta_j)$ by the formulation of Equation 3

Also, using the induction hypothesis, we can observe that the right hand side of 3 is increasing in α_j and decreasing in β_j because all terms there are increasing in α_j and decreasing in β_j . \square

This lemma enables us to prove that algorithm 3 is deciding a valid upper bound on the probability correctly guessing c examples out of c' guesses. To prove this, assume that the probability of such event is equal to τ' , Note that this means $\alpha_c + \beta_c = \frac{c'}{m} \tau'$. Also note that

$$\frac{\alpha_c}{\beta_c} \geq \frac{c}{c' - c} \quad (6)$$

therefore, we have

$$\alpha_c \geq \frac{c}{m} \tau' \quad (7)$$

and $\beta_c \leq \frac{c' - c}{m} \tau'$. Therefore, using Lemma 11 we have $\alpha_0 \geq r_{0,c}(\frac{c}{m} \tau', \frac{c' - c}{m} \tau')$ and $\beta_0 \geq h_{0,c}(\frac{c}{m} \tau', \frac{c' - c}{m} \tau')$.

Now we prove a lemma about the function $s_{i,j}(\tau) = h_{i,j}(\frac{c}{m} \tau, \frac{c' - c}{m} \tau) + r_{i,j}(\frac{c}{m} \tau, \frac{c' - c}{m} \tau)$.

Lemma 15. *the function $s_{i,j}(\tau) = h_{i,j}(\frac{c}{m} \tau, \frac{c' - c}{m} \tau) + r_{i,j}(\frac{c}{m} \tau, \frac{c' - c}{m} \tau)$ is increasing in τ for $i < j \leq c$.*

Proof. To prove this, we show that for all $i < j \leq c$ both $r_{i,j}(\frac{c}{m} \tau, \frac{c' - c}{m} \tau)$ and $h_{i,j}(\frac{c}{m} \tau, \frac{c' - c}{m} \tau)$ are increasing in τ . We prove this by induction on $j - i$. For $j - i = 1$, we have

$$h_{i,i+1}(\frac{c}{m} \tau, \frac{c' - c}{m} \tau) = (k - 1) \bar{f}^{-1}(\frac{c}{m} \tau).$$

We know that \bar{f}^{-1} is increasing, therefore $h_{i,i+1}(\frac{c}{m} \tau, \frac{c' - c}{m} \tau)$ is increasing in τ as well. For $r_{i,i+1}$ we have

$$r_{i,i+1}(\frac{c}{m} \tau, \frac{c' - c}{m} \tau) = \frac{c}{m} \tau + \frac{i}{c' - i} (h_{i,i+1}(\frac{c}{m} \tau, \frac{c' - c}{m} \tau) - \frac{c' - c}{m} \tau)$$

So we have

$$\begin{aligned} r_{i,i+1}(\frac{c}{m} \tau, \frac{c' - c}{m} \tau) &= \frac{c(c' - i) - i(c' - c)}{m(c' - i)} \tau + \frac{i}{c' - i} h_{i,i+1}(\frac{c}{m} \tau, \frac{c' - c}{m} \tau) \\ &= \frac{(c - i)c'}{m(c' - i)} \tau + \frac{i}{c' - i} h_{i,i+1}(\frac{c}{m} \tau, \frac{c' - c}{m} \tau). \end{aligned}$$

We already proved that $h_{i,i+1}(\frac{c}{m} \tau, \frac{c' - c}{m} \tau)$ is increasing in τ . We also have $\frac{(c - i)c'}{m(c' - i)} > 0$, since $i < c$. Therefore

$$r_{i,i+1}(\frac{c}{m} \tau, \frac{c' - c}{m} \tau)$$

is increasing in τ . So the base of induction is proved. Now we focus on $j - i > 1$. For $h_{i,j}$ we have

$$h_{i,j}(\frac{c}{m} \tau, \frac{c' - c}{m} \tau) = (k - 1) \bar{f}^{-1}(r_{i+1,j}(\frac{c}{m} \tau, \frac{c' - c}{m} \tau)).$$

By the induction hypothesis, we know that $r_{i+1,j}(\frac{c}{m}\tau, \frac{c'-c}{m}\tau)$ is increasing in τ , and we know that \bar{f}^{-1} is increasing, therefore, $h_{i,j}(\frac{c}{m}\tau, \frac{c'-c}{m}\tau)$ is increasing in τ .

For $r_{i,j}$, note that we rewrite it as follows

$$r_{i,j}(\alpha, \beta) = \alpha - \frac{j}{c' - j}\beta + \sum_{z=i}^{j-1} \lambda_z \cdot h_{z,j}(\alpha, \beta)$$

where $\lambda_z = (\frac{z+1}{c'-z-1} - \frac{z}{c'-z}) \geq 0$. Therefore, we have

$$\begin{aligned} r_{i,j}(\frac{c}{m}\tau, \frac{c'-c}{m}\tau) &= \tau(\frac{c}{m} - \frac{(c'-c)j}{m(c'-j)}) + \sum_{z=i}^{j-1} \lambda_z \cdot h_{z,j}(\frac{c}{m}\tau, \frac{c'-c}{m}\tau) \\ &= \tau \frac{c'(c-j)}{m(c'-j)} + \sum_{z=i}^{j-1} \lambda_z \cdot h_{z,j}(\frac{c}{m}\tau, \frac{c'-c}{m}\tau). \end{aligned}$$

Now we can verify that all terms in this equation are increasing in τ , following the induction hypothesis and the fact that $\lambda_z > 0$ and also $j \leq c$. \square

Now using this Lemma, we finish the proof. Note that we have $\alpha_0 + \beta_0 = \frac{c'}{m}$.

So assuming that $\tau' \geq \tau$, then we have

$$\frac{c'}{m} = \alpha_0 + \beta_0 \geq s_{0,c}(\tau') \geq s_{0,c}(\tau).$$

The last step of algorithm checks if $s_{0,c} \geq \frac{c'}{m}$ and it concludes that $\tau' \leq \tau$ if that's the case, because $s_{0,c}$ is increasing in τ . This means that the probability of having more than c guesses cannot be more than τ . \square

B ABLATION EXPERIMENTS

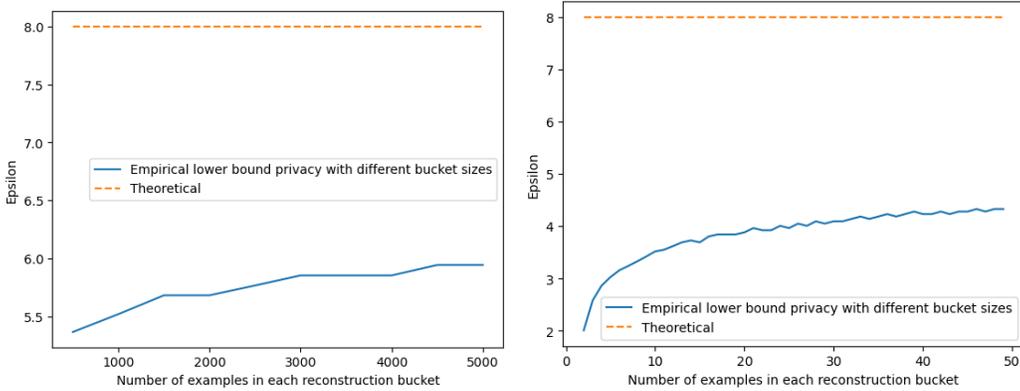


Figure 5: Effect of bucket size on the empirical lower bounds for reconstruction attack (Gaussian mechanism with standard deviation 0.6). Left: 10,000 canaries with bucket size up-to 5000. Right: 100 canaries with bucket-size up-to 50.

Reconstruction attacks: To show the effect of the bucket size (k) on the auditing performance, in Figure 5, we change the number of examples in the two different setups. In first setup we use 10,000 canaries and change the bucket size from 50 to 5000. In the other setup we only use 100 canaries and change the bucket-size from 3 to 50. Note that in these experiments, we do not use abstention and only consider adversaries that guess all examples.

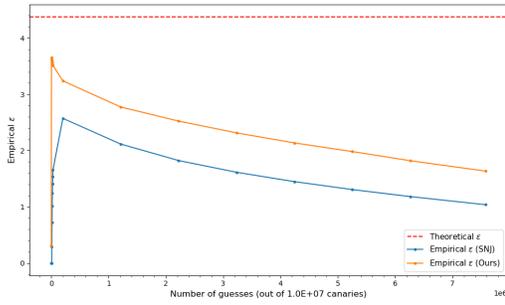


Figure 6: Effect of number of guesses (Gaussian mechanism with standard deviation 1.0)

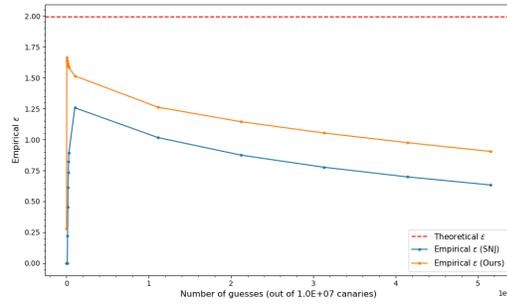


Figure 7: Effect of number of guesses (Gaussian mechanism with standard deviation 2.0)

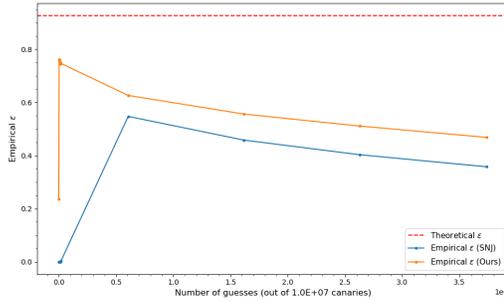


Figure 8: Effect of number of guesses (Gaussian mechanism with standard deviation 4.0)

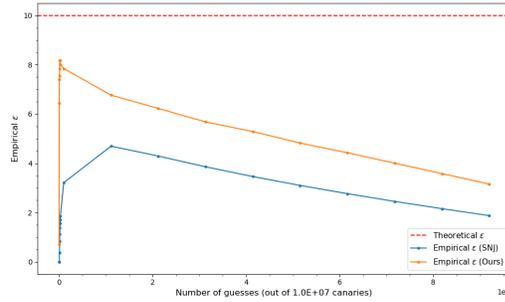


Figure 9: Effect of number of guesses (Gaussian mechanism with standard deviation 0.5)

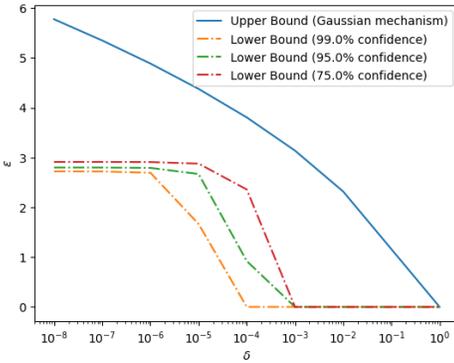


Figure 10: Idealized setting for different values of δ and confidence levels for bounds of Steinke et al. (2023).

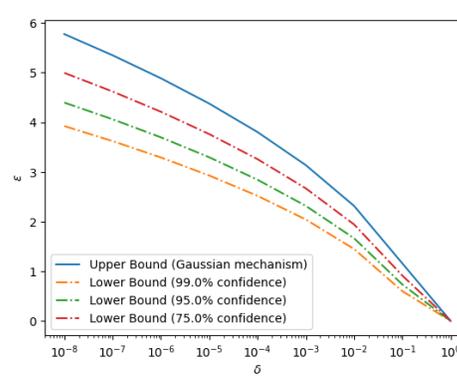


Figure 11: Idealized setting for different values of δ and confidence levels for our bounds.

Effect of number of guesses In Figures 6–9, we compare the theoretical upper bound, our lower bound and the bound of Steinke et al. lower bound with varying number of guesses. In total, we have $m = 10^7$ canaries. The number of correct guesses is determined by using Algorithm 4 (the idealized setting). Then we use our and Steinke et al. (2023)’s auditing with the resulting numbers and report the empirical ϵ . As we can see, both our and Steinke et al.’s auditing procedure achieve the best auditing performance for small number of guesses. This shows the importance of abstention in auditing.

A curious reader might wonder why the number of guesses has such a big impact on empirical privacy. Essentially, our analysis involves estimating how many correct guesses an adversary can make when given a certain number of attempts. We focus on specific percentiles of these distributions. The accuracy of our empirical privacy estimates can vary significantly based on how much the number

1188 of correct guesses fluctuates, which is influenced by how many guesses we allow the adversary to
1189 make. To explain further, consider a random variable representing the ratio of correct guesses (c) to
1190 total guesses (c'). If we reduce the number of guesses the variance of this ratio tends to decrease
1191 because the ratio approaches 1 (the adversary can make more correct guesses when we decrease c').
1192 Conversely, if we increase the number of guesses, the variance can also decrease because having more
1193 guesses generally leads to a more stable average, owing to the law of large numbers. This balance
1194 makes the number of guesses a crucial factor to optimize for the best estimation of empirical privacy.

1195
1196 **Varying δ and confidence levels:** We also examine the effect of δ on the obtained empirical ϵ .
1197 We fix the number of canaries to 10^5 and the number of guesses to 1, 500 and the number of correct
1198 guesses are set to 1, 429, suggested by the idealized setting. We use a Gaussian mechanism with
1199 standard deviation 1.0, we vary the value of δ and the confidence level to observe how they affect the
1200 results. Figures 10 and 11 shows the bound of Steinke et al. (2023) and our bound, respectively. Note
1201 that our lower bounds represent the true behavior of δ independent of the confidence level, in contrast
1202 to the bound of Steinke et al. (2023).

1203 C EXPERIMENTAL DETAILS

1205
1206 **Idealized setting:** In the idealized setting, we work with a toy version of the mechanism to calculate
1207 the *expected* number of correct guesses for the ideal adversary. For Gaussian mechanism, the ideal
1208 setting for an adversary is when we have a Gaussian mechanism that is used to calculate the sum
1209 of vectors. In this setting, each canary represents a unit vector that is orthogonal to all other canary
1210 vectors. Then, given the noisy sum, the adversary will calculate the likelihood of the canary being
1211 used in the sum, and then decides on the guesses based on these likelihoods. For the setting that
1212 the adversary has more than 2 guesses ($k > 2$), we use a slightly different idealized setting. In all
1213 settings, we run the attack 100 times and average the result to get the expected number of correct
1214 guesses. Algorithm 4 shows how we calculate the number of correct guesses in the idealized setting.

1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

Algorithm 4 Simulate the Number of Correct Guesses

```

import numpy as np
from scipy.special import softmax
from numpy.random import normal, binomial
def idealized_setting(target_noise, n_guesses, n_canaries, k):
    n_correct_vec = []
    if k==2:
        for _ in range(100):
            s_vector = binomial(1, 0.5, size=n_canaries) * 2 - 1
            noise = normal(0, 2*target_noise, n_canaries)

            noisy_s = s_vector + noise

            sorted_noisy_s = np.sort(noisy_s)

            threshold_c = sorted_noisy_s[-int(n_guesses)//2-1]
            n_correct = np.ceil(n_guesses*(s_vector[noisy_s >
                ↪ threshold_c] == 1).mean())

            n_correct_vec.append(n_correct)
    else:
        for _ in range(100):
            s_recon_vec = np.random.randint(0, k, n_canaries)

            s_vec_recn_ohc = np.eye(k)[s_recon_vec]
            s_recon_noisy_vec_ohc = s_vec_recn_ohc + normal(0,
                ↪ np.sqrt(2)*target_noise, s_vec_recn_ohc.shape)

            idx_max = np.argmax(s_recon_noisy_vec_ohc, axis=1)

            buckets =
                ↪ softmax(s_recon_noisy_vec_ohc/(2*target_noise**2),
                ↪ axis=1)[np.arange(s_recon_noisy_vec_ohc.shape[0]),
                ↪ idx_max]
            sorted_buckets = np.sort(buckets)
            bucket_c_thr = sorted_buckets[-int(n_guesses)]

            n_correct_rec = np.ceil(
                n_guesses*(s_recon_vec[buckets > bucket_c_thr] ==
                ↪ s_recon_noisy_vec_ohc[buckets >
                ↪ bucket_c_thr].argmax(1)).mean()
            )
            n_correct_vec.append(n_correct_rec)

    return int(np.array(n_correct_vec).mean(0))

```

1296 AUDITING CODE

1297

1298

1299 Here we include the code to compute empirical epsilon.

1300

1301

1302

1303

1304

1305

1306

1307

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324

1325

1326

1327

1328

1329

1330

1331

1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

```

from scipy.stats import norm
import numpy as np

# Calculate h and r recursively (no abstentions)
def rh(inverse_blow_up_function, alpha, beta, j, m, k=2):
    # Initialize lists to store h and r values
    h = [0 for _ in range(j + 1)]
    r = [0 for _ in range(j + 1)]
    # Set initial values for h and r
    h[j] = beta
    r[j] = alpha
    # Iterate from j-1 to 0
    for i in range(j - 1, -1, -1):
        # Calculate h[i] using the maximum of h[i+1] and a
        ↪ scaled inverse blow-up function
        h[i] = max(h[i + 1], (k - 1) *
        ↪ inverse_blow_up_function(r[i + 1]))
        # Update r[i] based on the difference between h[i] and
        ↪ h[i+1]
        r[i] = r[i + 1] + (i / (m - i)) * (h[i] - h[i + 1])
    # Return the lists of h and r values
    return (r, h)

# Audit function without abstention
def audit_rh(inverse_blow_up_function, m, c, threshold=0.05,
    ↪ k=2):
    # Calculate alpha and beta values
    alpha = threshold * c / m
    beta = threshold * (m - c) / m
    # Call the rh function to get the lists of h and r values
    r, h = rh(inverse_blow_up_function, alpha, beta, c, m, k)
    # Check if the differential privacy condition is satisfied
    if r[0] + h[0] > 1.0:
        return False
    else:
        return True

# Calculate h and r recursively (with abstentions)
def rh_with_cap(inverse_blow_up_function, alpha, beta, j,
    ↪ m, c_cap, k=2):
    h=[0 for i in range(j+1)]
    r=[0 for i in range(j+1)]
    h[j]= beta
    r[j]= alpha
    for i in range(j-1,-1,-1):
        ↪ h[i]=max(h[i+1],(k-1)*inverse_blow_up_function(r[i+1]))
        r[i]= r[i+1] + (i/(c_cap-i))*(h[i] - h[i+1])

    return (r,h)

# Audit function with abstentions
def audit_rh_with_cap(inverse_blow_up_function, m, c, c_cap,
    ↪ threshold=0.05, k=2):
    threshold=threshold*c_cap/m

```

```

1350     alpha=(threshold*c/c_cap)
1351     beta=threshold*(c_cap-c)/c_cap
1352     r,h=rh_with_cap(inverse_blow_up_function, alpha, beta, c,
1353     ↪ m, c_cap, k)
1354
1355     if r[0]+h[0]>c_cap/m:
1356         return False
1357     else:
1358         return True
1359
1360     # Calculate the blow-up function for Gaussian noise
1361     def gaussianDP_blow_up_function(noise):
1362         def blow_up_function(x):
1363             # Calculate the threshold value
1364             threshold = norm.ppf(x)
1365             # Calculate the blown-up threshold value
1366             blown_up_threshold = threshold + 1 / noise
1367             # Return the CDF of the blown-up threshold value
1368             return norm.cdf(blown_up_threshold)
1369         return blow_up_function
1370
1371     # Calculate the inverse blow-up function for Gaussian noise
1372     def gaussianDP_blow_up_inverse(noise):
1373         def blow_up_inverse_function(x):
1374             # Calculate the threshold value
1375             threshold = norm.ppf(x)
1376             # Calculate the blown-up threshold value
1377             blown_up_threshold = threshold - 1 / noise
1378             # Return the CDF of the blown-up threshold value
1379             return norm.cdf(blown_up_threshold)
1380         return blow_up_inverse_function
1381
1382     # Define a function to calculate delta for Gaussian noise
1383     def calculate_delta_gaussian(noise, epsilon):
1384         # Calculate delta using the formula
1385         delta = norm.cdf(-epsilon * noise + 1 / (2 * noise)) -
1386         ↪ np.exp(epsilon) * norm.cdf(-epsilon * noise - 1 / (2 *
1387         ↪ noise))
1388         return delta
1389
1390     # Define a function to calculate epsilon for Gaussian noise
1391     def calculate_epsilon_gaussian(noise, delta):
1392         # Set initial bounds for epsilon
1393         epsilon_upper = 100
1394         epsilon_lower = 0
1395         # Perform binary search to find epsilon
1396         while epsilon_upper - epsilon_lower > 0.001:
1397             epsilon_middle = (epsilon_upper + epsilon_lower) / 2
1398             if calculate_delta_gaussian(noise, epsilon_middle) >
1399             ↪ delta:
1400                 epsilon_lower = epsilon_middle
1401             else:
1402                 epsilon_upper = epsilon_middle
1403         # Return the upper bound of epsilon
1404         return epsilon_upper
1405
1406     # Get the empirical epsilon value
1407     def get_gaussian_emp_eps_ours(candidate_noises,
1408     ↪ inverse_blow_up_functions, m, c, threshold, delta, k=2):

```

```
1404     # Initialize the empirical privacy index
1405     empirical_privacy_index = 0
1406     # Iterate through candidate noises until the privacy
1407     ↪ condition fails
1408     while
1409     ↪ audit_rh(inverse_blow_up_functions[empirical_privacy_index],
1410     ↪ m, c, threshold=0.05, k=k):
1411         empirical_privacy_index += 1
1412     # Get the empirical noise and calculate the empirical
1413     ↪ epsilon
1414     empirical_noise =
1415     ↪ candidate_noises[empirical_privacy_index]
1416     empirical_eps =
1417     ↪ calculate_epsilon_gaussian(empirical_noise,
1418     ↪ delta=delta)
1419     # Return the empirical epsilon
1420     return empirical_eps
1421
1422     # Set target noise and generate candidate noises
1423     target_noise = 0.6
1424     candidate_noises=[target_noise+ i*0.01 for i in range(1000)]
1425     inverse_blow_up_functions=[gaussianDP_blow_up_inverse(noise)
1426     ↪ for noise in candidate_noises]
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
```