
One-shot Empirical Privacy Estimation for Federated Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Privacy estimation techniques for differentially private (DP) algorithms are useful
2 for comparing against analytical bounds, or to empirically measure privacy loss in
3 settings where known analytical bounds are not tight. However, existing privacy
4 auditing techniques usually make strong assumptions on the adversary (e.g., knowl-
5 edge of intermediate model iterates or the training data distribution), are tailored to
6 specific tasks, model architectures, or DP algorithm, and/or require retraining the
7 model many times (typically on the order of thousands). These shortcomings make
8 deploying such techniques at scale difficult in practice, especially in federated
9 settings where model training can take days or weeks. In this work, we present a
10 novel “one-shot” approach that can systematically address these challenges, allow-
11 ing efficient auditing or estimation of the privacy loss of a model during the same,
12 single training run used to fit model parameters, and without requiring any *a priori*
13 knowledge about the model architecture, task, or DP training algorithm. We show
14 that our method provides provably correct estimates for the privacy loss under the
15 Gaussian mechanism, and we demonstrate its performance on well-established FL
16 benchmark datasets under several adversarial threat models.

17 1 Introduction

18 Federated learning (FL) [McMahan et al., 2017, Kairouz et al., 2021b] is a paradigm for training
19 machine learning models on decentralized data. At each round, selected clients contribute model
20 updates to be aggregated by a server, without ever communicating their raw data. FL incorporates
21 *data minimization* principles to reduce the risk of compromising anyone’s data: each user’s data
22 never leaves their device, the update that is transmitted contains only information necessary to update
23 the model, the update is encrypted in transit, and the update exists only ephemerally before being
24 combined with other clients’ updates and then incorporated into the model [Bonawitz et al., 2022].
25 Technologies such as secure aggregation [Bonawitz et al., 2017, Bell et al., 2020] can be applied to
26 ensure that even the central server cannot inspect individual updates, but only their aggregate.

27 However, these data minimization approaches cannot rule out the possibility that an attacker might
28 learn some private information from the training data by directly interrogating the final model [Carlini
29 et al., 2021, Balle et al., 2022, Haim et al., 2022]. To protect against this, *data anonymization* for the
30 model is required. FL can be augmented to satisfy user-level differential privacy [Dwork and Roth,
31 2014, Abadi et al., 2016, McMahan et al., 2018], the gold-standard for data anonymization. DP can
32 guarantee each user that a powerful attacker – one who knows all other users’ data, all details about
33 the algorithm (other than the values of the noise added for DP), and every intermediate model update
34 – still cannot confidently infer the presence of that user in the population, or anything about their data.
35 This guarantee is typically quantified by the parameter ϵ , with lower values corresponding to higher
36 privacy (less confidence for the attacker).

37 DP is often complemented by *empirical privacy estimation* techniques, such as membership inference
38 attacks [Shokri et al., 2017, Yeom et al., 2018, Carlini et al., 2022], which measure the success of an
39 adversary at distinguishing whether a particular record was part of training or not.¹ Such methods
40 have been used to audit the implementations of DP mechanisms or claims about models trained with
41 DP [Jagielski et al., 2020, Nasr et al., 2021, Zanella-Béguelin et al., 2022, Lu et al., 2022]. They
42 are also useful for estimating the privacy loss in cases where a tight analytical upper bound on ϵ is
43 unknown, for example when clients are constrained to participate in at most some number of rounds,
44 or when the adversary does not see the full trace of model iterates. However, existing privacy auditing
45 techniques suffer from several major shortcomings. First, they require retraining the model many
46 times (typically in the thousands) to provide reliable estimates of DP’s ϵ [Jagielski et al., 2020, Nasr
47 et al., 2021]. Second, they often rely on knowledge of the model architecture and/or the underlying
48 dataset (or at least a similar, proxy dataset) for mounting the attack. For example, a common approach
49 is to craft a “canary” training example on which the membership is being tested, which typically
50 requires an adversary to have access to the underlying dataset and knowledge of the domain and
51 model architecture. Finally, such techniques typically grant the adversary unrealistic power, for
52 example (and in particular) the ability to inspect all model iterates during training [Maddock et al.,
53 2022], something which may or may not be reasonable depending on the system release model.

54 Such assumptions are particularly difficult to satisfy in FL due to the following considerations:

- 55 • **Minimal access to the dataset, or even to proxy data.** A primary motivating feature of FL is is
56 that it can make use of on-device data without (any) centralized data collection. In many tasks,
57 on-device data is more representative of real-world user behavior than any available proxy data.
- 58 • **Infeasibility of training many times, or even more than one time.** FL training can take days or
59 weeks, and expends resources on client devices. To minimize auditing time and client resource
60 usage, an ideal auditing technique should produce an estimate of privacy during the same, single
61 training run used to optimize model parameters, and without significant overhead from crafting
62 examples or computing additional “fake” training rounds.
- 63 • **Lack of task, domain, and model architecture knowledge.** A scalable production FL platform
64 is expected to cater to the needs of many diverse ML applications, from speech to image to
65 language modeling tasks. Therefore, using techniques that require specific knowledge of the task
66 and/or model architecture makes it hard to deploy those techniques at scale in production settings.

67 In this paper, we design an auditing technique tailored for FL usage with those considerations in mind.
68 We empirically estimate ϵ efficiently under user-level DP federated learning by measuring the training
69 algorithm’s tendency to memorize arbitrary clients’ updates. Our main insight is to insert multiple
70 canary clients in the federated learning protocol with independent random model updates, and design
71 a test statistic based on cosine angles of each canary update with the final model to test participation
72 of a certain user in the protocol. The intuition behind the approach comes from the elementary
73 result that in a high-dimensional space, isotropically sampled vectors are nearly orthogonal with high
74 probability. So we can think of each canary as estimating the algorithm’s tendency to memorize along
75 a dimension of variance that is independent of the true model updates, and of the other canaries.

76 Our method has several favorable properties. It can be applied during the same, single training
77 run which is used to train the federated model parameters, and therefore does not incur additional
78 performance overhead. Although it does inject some extra noise into the training process, the effect on
79 model quality is negligible, provided model dimensionality and number of clients are reasonably sized.
80 We show that in the tractable case of a single application of the Gaussian mechanism, our method
81 provably recovers the true, analytical ϵ in the limit of high dimensionality. We evaluate privacy loss
82 for several adversarial models of interest, for which existing analytical bounds are not tight. In the
83 case when all intermediate updates are observed and the noise is low, our method produces high values
84 of ϵ , indicating that an attacker could successfully mount a membership inference attack. However,
85 in the common and important case that only the final trained model is released, our ϵ estimate is far
86 lower, suggesting that adding a modest amount of noise is sufficient to prevent leakage, as has been
87 observed by practitioners. Our method can also be used to explore how leakage changes as aspects of
88 the training protocol change, for which no tight theoretical analysis is known, for example if we limit
89 client participation. The method we propose is model and dataset agnostic, so it can be easily applied
90 without change to any federated learning task.

¹Some prior work only applies to example-level DP, in which *records* correspond to examples, as opposed to user-level, in which *records* are users. We will describe our approach in terms of user-level DP, but it can be trivially modified to provide example-level DP.

91 **2 Background and related work**

92 **Differential privacy.** Differential privacy (DP) [Dwork et al., 2006, Dwork and Roth, 2014] is
 93 a rigorous notion of privacy that an algorithm can satisfy. DP algorithms for training ML models
 94 include DP-SGD [Abadi et al., 2016], DP-FTRL [Kairouz et al., 2021a], and DP matrix factoriza-
 95 tion [Denissoff et al., 2022, Choquette-Choo et al., 2022]. Informally, DP guarantees that a powerful
 96 attacker observing the output of the algorithm A trained on one of two *adjacent* datasets (differing by
 97 addition or removal of one record), D or D' , cannot confidently distinguish the two cases, which is
 98 quantified by the privacy parameter ϵ .

Definition 2.1. User-level differential privacy. The training algorithm $A : \mathcal{D} \rightarrow \mathcal{R}$ is user-level
 (ϵ, δ) differentially private if for all pairs of datasets D and D' from \mathcal{D} that differ only by addition or
 removal of the data of one user and all output regions $R \subseteq \mathcal{R}$:

$$\Pr[A(D) \in R] \leq e^\epsilon \Pr[A(D') \in R] + \delta.$$

99 DP can be interpreted as a hypothesis test with the null hypothesis that A was trained on D and the
 100 alternative hypothesis that A was trained on D' . False positives (type-I errors) occur when the null
 101 hypothesis is true, but is rejected, while false negatives (type-II errors) occur when the alternative
 102 hypothesis is true, but is rejected. Kairouz et al. [2015] characterized (ϵ, δ)-DP in terms of the
 103 false positive rate (FPR) and false negative rate (FNR) achievable by an acceptance region. This
 104 characterization enables estimating the privacy parameter as:

$$\hat{\epsilon} = \max\left\{\log \frac{1 - \delta - \text{FPR}}{\text{FNR}}, \log \frac{1 - \delta - \text{FNR}}{\text{FPR}}\right\}. \quad (1)$$

105 We review and compare with related work in Appendix F.

106 **3 One-shot privacy estimation for the Gaussian mechanism**

107 As a warm-up, we start by considering the problem of estimating the privacy of the Gaussian
 108 mechanism, the fundamental building block of DP-SGD and DP-FedAvg. To be precise, given
 109 $D = (x_1, \dots, x_n)$, with $\|x_i\| \leq 1$ for all $i \in [n]$, the output of the Gaussian vector sum query is
 110 $A(D) = \bar{x} + \sigma Z$, where $\bar{x} = \sum_i x_i$ and $Z \sim \mathcal{N}(0, I)$. Without loss of generality, we can consider
 111 a neighboring dataset D' with an additional vector x with $\|x\| \leq 1$. Thus, $A(D) \sim \mathcal{N}(\bar{x}, \sigma^2 I)$
 112 and $A(D') \sim \mathcal{N}(\bar{x} + x, \sigma^2 I)$. For the purpose of computing the DP guarantees, this mechanism is
 113 equivalent to analyzing $A(D) \sim \mathcal{N}(0, \sigma^2)$ and $A(D') \sim \mathcal{N}(1, \sigma^2)$ due to spherical symmetry.

114 The naive approach for estimating the ϵ of an implementation of the Gaussian mechanism would
 115 run it many times (say 1000 times), with half of the runs on D and the other half on D' . Then the
 116 outputs of these runs are shuffled and given to an ‘‘attacker’’ who attempts to determine for each
 117 output whether it was computed from D or D' . Finally, the performance of the attacker is quantified
 118 and Eq. (1) is used to obtain an estimate of the mechanism’s ϵ at a target δ .

119 We now present an approach for estimating ϵ by running the mechanism *only once*. The basic idea
 120 behind our approach is to augment the original dataset with k canary vectors c_1, \dots, c_k , sampled i.i.d
 121 uniformly at random from the unit sphere, obtaining $D = (x_1, \dots, x_n, c_1, \dots, c_k)$. We consider
 122 k neighboring datasets, each excluding one of the canaries, i.e., $D'_i = D \setminus \{c_i\}$ for $i \in [k]$. We
 123 run the Gaussian mechanism once on D and use its output to compute k test statistics $\{g_i\}_{i \in [k]}$, the
 124 cosine of the angles between the output and each one of the k canary vectors. We use these k cosines
 125 to estimate the distribution of test statistic on D by computing the sample mean $\hat{\mu} = \frac{1}{k} \sum_{j=1}^k g_j$
 126 and sample variance $\hat{\sigma}^2 = \frac{1}{k} \sum_{j=1}^k (g_j - \hat{\mu})^2$ and fitting a Gaussian $\mathcal{N}(\hat{\mu}, \hat{\sigma}^2)$. To estimate the
 127 distribution of the test statistic on D' , we need to run the mechanism on each D'_i and compute the
 128 cosine of the angle between the output vector and c_i . This is where our choice of (i) independent
 129 isotropically distributed canaries and (ii) cosine angles as our test statistic are particularly useful.
 130 The distribution of the cosine of the angle between an isotropically distributed unobserved canary
 131 and the mechanism output (or any independent vector) can be described in a closed form; there is no
 132 need to approximate this distribution with samples. We will show in Theorems 3.1 and 3.2 that this
 133 distribution can be well approximated by $\mathcal{N}(0, 1/d)$. Now that we have the distribution of the test
 134 statistic on D and D' , we estimate the ϵ of the mechanism using the method given in Appendix B

135 which allows us to compute the ε when the null and alternate hypotheses are two arbitrary Gaussians.
 136 Our approach is summarized in Algorithm 1.

Algorithm 1 One-shot privacy estimation for Gaussian mechanism.

1: **Input:** Vectors x_1, \dots, x_n with $\|x_i\| \leq 1$,
 DP noise variance σ^2 , and target δ
 2: $\rho \leftarrow \sum_{i \in [n]} x_i$
 3: **for** $j \in [k]$ **do**
 4: Draw random $c_j \in \mathbb{S}^{d-1}$ unit sphere
 5: $\rho \leftarrow \rho + c_j$
 6: **Release** $\rho \leftarrow \rho + \mathcal{N}(0, \sigma^2 I)$
 7: **for** $j \in [k]$ **do**
 8: $g_j \leftarrow \langle c_j, \rho \rangle / \|\rho\|$
 9: $\hat{\mu}, \hat{\sigma} \leftarrow \mathbf{mean}(\{g_j\}), \mathbf{std}(\{g_j\})$
 10: $\hat{\varepsilon} \leftarrow \varepsilon(\mathcal{N}(0, 1/d) \parallel \mathcal{N}(\hat{\mu}, \hat{\sigma}^2); \delta)$

137 We argue that the approach given in Algorithm 1 gives an estimate of ε that approaches the exact
 138 value when d is high. To do so, we will prove (Theorems 3.1 and 3.2) that distribution of the test
 139 statistic on D' is indeed well approximated by $\mathcal{N}(0, 1/d)$, and (Theorem 3.3) that $\sqrt{d}\hat{\mu} \xrightarrow{p} 1/\sigma$ and
 140 $d\hat{\sigma}^2 \xrightarrow{p} 1$ as $d \rightarrow \infty$, i.e., the distribution of test statistic on D is $\mathcal{N}(\frac{1}{\sigma\sqrt{d}}, \frac{1}{d})$, asymptotically. Since
 141 these two distributions are just a scaling of $A(D) \sim \mathcal{N}(0, \sigma^2)$ and $A(D') \sim \mathcal{N}(1, \sigma^2)$ by a factor of
 142 $\frac{1}{\sigma\sqrt{d}}$, the ε is the same, which proves our claim.² (All proofs in Appendix A.)

143 **Theorem 3.1.** *Let S be sampled uniformly from the unit sphere in \mathbb{R}^d , and let $\tau = \langle S, v \rangle / \|v\| \in$
 144 $[-1, 1]$ be the cosine similarity between S and some arbitrary independent nonzero vector v . Then,
 145 the probability density function of τ is*

$$f_d(\tau) = \frac{\Gamma(\frac{d}{2})}{\Gamma(\frac{d-1}{2})\sqrt{\pi}}(1 - \tau^2)^{\frac{d-3}{2}}.$$

146 **Theorem 3.2.** *In the setting of Theorem 3.1, we have that $\tau\sqrt{d}$ converges in distribution to $\mathcal{N}(0, 1)$
 147 as $d \rightarrow \infty$, i.e., $\lim_{d \rightarrow \infty} \mathbb{P}(\tau \leq \lambda/\sqrt{d}) = \mathbb{P}_{Z \sim \mathcal{N}(0,1)}(Z \leq \lambda)$.*

148 **Theorem 3.3.** *For $d \in \mathbb{N}$, let $k = o(d)$, but $k = \omega(1)$. For $i = 1 \dots k$, let c_i sampled i.i.d.
 149 from the unit sphere in d dimensions. Let $Z \sim \mathcal{N}(0; I_d)$. Let $\sigma > 0$, and define the mechanism
 150 result $\rho = \sum_{j=1}^k c_j + \sigma Z$, and the cosine values $g_j = \frac{\langle c_j, \rho \rangle}{\|\rho\|}$. Write the empirical mean of the
 151 cosines $\hat{\mu} = \frac{1}{k} \sum_{j=1}^k g_j$, and the empirical variance $\hat{\sigma}^2 = \frac{1}{k} \sum_{j=1}^k (g_j - \hat{\mu})^2$. Then as $d \rightarrow \infty$,
 152 $\sqrt{d}\hat{\mu} \xrightarrow{p} 1/\sigma$ and $d\hat{\sigma}^2 \xrightarrow{p} 1$.*

153 We note that running the algorithm with moderate values of d and k already yields a close approxi-
 154 mation. We simulated the case where $d = 10^6$, $k = 10^3$, and the results are shown in Table 1. The
 155 estimated ε_{est} is very close to the true value of ε , with small standard deviation, which demonstrates
 156 that the random canary method provides tight estimates for the Gaussian mechanism.

σ	analytical ε	ε_{est}
4.22	1.0	0.972 ± 0.148
1.54	3.0	3.04 ± 0.137
0.541	10.0	9.98 ± 0.190

Table 1: One-shot auditing of the Gaussian mechanism with $d = 10^6$, $k = 10^3$, and $\delta = 10^{-6}$. For each value of ε , we set σ using the optimal calibration of Balle and Wang [2018], and then use the random canary method to output the estimate ε_{est} . Shown is the mean and std ε_{est} over 50 simulations.

157 4 One-shot privacy estimation for FL with random canaries

158 We now extend this idea to DP Federated Averaging to estimate the privacy of releasing the final
 159 model parameters in one shot, during model training. We propose adding k canary clients to the

²Our Theorem 3.3 as written does not include the data vectors x_i in the sum ρ . It can be extended to do so if we also assume that $n = o(d)$.

160 training population who participate exactly as real users do. Each canary client generates a random
161 model update sampled from the unit sphere, which it returns at every round in which it participates,
162 scaled to have norm equal to the clipping norm for the round. After training, we collect the set
163 of canary/final-model cosines, fit them to a Gaussian, and compare them to the null hypothesis
164 distribution $\mathcal{N}(0, 1/d)$ just as we did for the basic Gaussian mechanism. The procedure is described
165 in Algorithm 2. (Algorithms 2 and 3 can be found in Appendix C.)

166 FL is an optimization procedure in which each model iterate is a linear combination of all updates
167 received thus far, plus Gaussian noise. Our threat model allows the attacker to control the updates of
168 a client, and the ability to inspect the final model. We argue that it is a powerful (perhaps optimal,
169 under some assumptions) strategy to return a very large update that is essentially orthogonal to all
170 other updates, and then measure the dot product (or cosine) to the final model. Here we use the fact
171 that randomly sampled canary updates are nearly orthogonal to all the true client updates and also
172 to each other. Unlike many works that only produce correct estimates when clients are sampled
173 uniformly and independently at each round, our method makes no assumptions on the pattern of client
174 participation. Clients may be sampled uniformly at each round, shuffled and processed in batches,
175 or even participate according to the difficult-to-characterize *de facto* pattern of participation of real
176 users in a production system. Our only assumption is that canaries can be inserted according to the
177 same distribution that real clients are. In production settings, a simple and effective strategy would
178 be to designate a small fraction of real clients to have their model updates replaced with the canary
179 update whenever they participate. If the participation pattern is such that memorization is easier, for
180 whatever reason, the distribution of canary/final-model cosines will have a higher mean, leading to
181 higher ε estimates.

182 We stress that our empirical ε_{est} estimate should not be construed as a formal bound on the worst-case
183 privacy leakage. Rather, a low value of ε_{est} can be taken as evidence that an adversary implementing
184 this particular, powerful attack will have a hard time inferring the presence of any given user upon
185 observing the final model. If we suppose that the attack is strong, or even optimal, then we can infer
186 that *any* attacker will not be able to perform MI successfully, and therefore our ε_{est} is a justifiable
187 metric of the true privacy when the final model is released. Investigating conditions under which this
188 could be proven would be a valuable direction for future work.

189 Aside from quantifying the privacy of releasing only the final model, our method allows us to explore
190 how privacy properties are affected by varying aspects of training for which we have no tight formal
191 analysis. As an important example (which we explore in experiments) we consider how the estimate
192 changes if clients are constrained to participate a fixed number of times.

193 We also propose a simple extension to our method that allows us to estimate ε under the threat model
194 where all model updates are observed. We use as the test statistic the *maximum over rounds* of the
195 angle between the canary and the model delta at that round. Unfortunately in this case we can no
196 longer express in closed form the distribution of max-over-rounds cosine of an canary that did not
197 participate in training, because it depends on the trajectory of partially trained models, which is task
198 and model specific. Our solution is to sample a set of unobserved canaries that are never included in
199 model updates, but we still keep track of their cosines with each model delta and finally take the max.
200 We approximate both the distributions of observed and unobserved maximum canary/model-delta
201 cosines using Gaussian distributions and compute the optimal ε . The pseudocode for this modified
202 procedure is provided in Algorithm 3. We will see that this method provides estimates of ε close to
203 the analytical bounds under moderate amounts of noise, providing evidence that our attack is strong.

204 5 Experiments

205 In this section we present the results of experiments estimating the privacy leakage while training
206 a model on a large-scale public federated learning dataset: the stackoverflow word prediction
207 data/model of Reddi et al. [2020]. The model is a word-based LSTM with 4.1M parameters. We train
208 the model for 2048 rounds with 167 clients per round, where each of the $m=34$ k clients participates
209 in exactly one round, amounting to a single epoch over the data. We use the adaptive clipping method
210 of Andrew et al. [2021]. With preliminary manual tuning, we selected a client learning rate of 1.0
211 and server learning rate of 0.56 for all experiments because the choice gives good performance over a
212 range of levels of DP noise. We always use 1k canaries for each set of cosines; experiments with
213 intermediate iterates use 1k observed and 1k unobserved canaries. We fix $\delta = m^{-1.1}$. We consider

Noise	analytical ϵ	ϵ_{lo} -all	ϵ_{est} -all	ϵ_{lo} -final	ϵ_{est} -final
0	∞	6.240	45800	2.88	4.60
0.0496	300	6.238	382	1.11	1.97
0.0986	100	5.05	89.4	0.688	1.18
0.2317	30	0.407	2.693	0.311	0.569

Table 2: Comparing ϵ estimates using all model deltas vs. using the final model only. ϵ_{lo} is the empirical 95% lower bound from our modified Jagielski et al. [2020] method. For moderate noise, ϵ_{est} -all is in the ballpark of the analytical ϵ , providing evidence that the attack is strong and therefore the ϵ estimates are reliable. On the other hand, ϵ_{est} -final is far lower, indicating that when the final model is observed, privacy is better.

214 noise multipliers³ in the range 0.0496 to 0.2317, corresponding to analytical ϵ estimates from 300
215 down to 30.⁴ We also include experiments with clipping only (noise multiplier is 0). We note that
216 across the range of noise multipliers, the participation of 1k canaries had no significant impact on
217 model accuracy – at most causing a 0.1% relative decrease.

218 We also report a high-probability lower bound on ϵ that comes from applying a modified version of
219 the method of Jagielski et al. [2020] to the set of cosines. That work uses Clopper-Pearson upper
220 bounds on the achievable FPR and FNR of a thresholding classifier to derive a bound on ϵ . We make
221 two changes: following Zanella-Béguelin et al. [2022], we use the tighter and more centered Jeffreys
222 confidence interval for the upper bound on FNR at some threshold a , and we use the exact CDF of
223 the null distribution for the FPR as described in Section 4. We refer to this lower bound as ϵ_{lo} . We
224 set $\alpha = 0.05$ to get a 95%-confidence bound.

225 We first consider the case where the intermediate updates are released as described in Algorithm 3.
226 The middle columns of Table 2 shows the results of these experiments over a range of noise multipliers.
227 For the lower noise multipliers, our method easily separates the cosines of observed vs. unobserved
228 canaries, producing very high estimates ϵ_{est} -all, which are much higher than lower bounds ϵ_{lo} -all
229 estimated by previous work. This confirms our intuition that intermediate model updates give the
230 adversary significant power to detect the presence of individuals in the data. It also provides evidence
231 that the canary cosine attack is strong, increasing our confidence that the ϵ estimates assuming a
232 weakened adversary that observes only the final model is not a severe underestimate.

233 The rightmost columns of Table 2 show the results of restricting the adversary to observe only the final
234 model, as described in Algorithm 2. Now ϵ_{est} is significantly smaller than when the adversary has
235 access to all intermediary updates. With clipping only, our estimate is 4.60, which is still essentially
236 vacuous from a rigorous privacy perspective.⁵ But with even a small amount of noise, we approach
237 the high-privacy regime of $\epsilon \sim 1$, confirming observations of practitioners that a small amount of
238 noise is sufficient to prevent most memorization.

239 In Appendix D we provide analogous experiments on the federated EMNIST dataset, with similar
240 results. In Appendix E we present further experiments to highlight the ability of our method to
241 estimate privacy when we vary aspects of the training algorithm that may reasonably be expected to
242 change privacy properties, but for which no tight analysis has been obtained.

243 6 Conclusion

244 We have introduced a novel method for empirically estimating the privacy loss during training of a
245 model with DP-FedAvg. For natural production-sized problems (millions of parameters, hundreds of
246 thousands of clients), it produces reasonable privacy estimates during the same single training run
247 used to estimate model parameters, without degrading the utility of the model, and does not require
248 any prior knowledge of the task, data or model.

³The noise multiplier is the ratio of the noise to the clip norm. When adaptive clipping is used, the clip norm varies across rounds, and the noise scales proportionally.

⁴Since each user participates once, we bound ϵ as the unamplified Gaussian mechanism applied once with no composition.

⁵An ϵ of 5 means that an attacker can go from a small suspicion that a user participated (say, 10%) to a very high degree of certainty (94%) [Desfontaines, 2018].

249 **References**

- 250 M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang. Deep
251 learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC conference on*
252 *computer and communications security*, pages 308–318, 2016.
- 253 G. Andrew, O. Thakkar, B. McMahan, and S. Ramaswamy. Differentially private learning with
254 adaptive clipping. *Advances in Neural Information Processing Systems*, 34:17455–17466, 2021.
- 255 B. Balle and Y.-X. Wang. Improving the gaussian mechanism for differential privacy: Analytical
256 calibration and optimal denoising. In *International Conference on Machine Learning*, pages
257 394–403. PMLR, 2018.
- 258 B. Balle, G. Cherubin, and J. Hayes. Reconstructing training data with informed adversaries. In *2022*
259 *IEEE Symposium on Security and Privacy (SP)*, pages 1138–1156. IEEE, 2022.
- 260 J. H. Bell, K. A. Bonawitz, A. Gascón, T. Lepoint, and M. Raykova. Secure single-server aggregation
261 with (poly) logarithmic overhead. In *Proceedings of the 2020 ACM SIGSAC Conference on*
262 *Computer and Communications Security*, pages 1253–1269, 2020.
- 263 K. Bonawitz, V. Ivanov, B. Kreuter, A. Marcedone, H. B. McMahan, S. Patel, D. Ramage, A. Segal,
264 and K. Seth. Practical secure aggregation for privacy-preserving machine learning. In *proceedings*
265 *of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, pages 1175–
266 1191, 2017.
- 267 K. Bonawitz, P. Kairouz, B. McMahan, and D. Ramage. Federated learning and privacy. *Commun.*
268 *ACM*, 65(4):90–97, mar 2022. ISSN 0001-0782. doi: 10.1145/3500240. URL [https://doi.](https://doi.org/10.1145/3500240)
269 [org/10.1145/3500240](https://doi.org/10.1145/3500240).
- 270 N. Carlini, F. Tramer, E. Wallace, M. Jagielski, A. Herbert-Voss, K. Lee, A. Roberts, T. Brown,
271 D. Song, U. Erlingsson, A. Oprea, and C. Raffel. Extracting training data from large language
272 models. In *30th USENIX Security Symposium (USENIX Security 2021)*, 2021.
- 273 N. Carlini, S. Chien, M. Nasr, S. Song, A. Terzis, and F. Tramer. Membership inference attacks
274 from first principles. In *IEEE Symposium on Security and Privacy (SP)*, pages 1519–1519, Los
275 Alamitos, CA, USA, May 2022. IEEE Computer Society. doi: 10.1109/SP46214.2022.00090. URL
276 <https://doi.ieeeecomputersociety.org/10.1109/SP46214.2022.00090>.
- 277 C. A. Choquette-Choo, H. B. McMahan, K. Rush, and A. Thakurta. Multi-epoch matrix factorization
278 mechanisms for private machine learning. *arXiv preprint arXiv:2211.06530*, 2022.
- 279 S. Denissov, H. B. McMahan, J. K. Rush, A. Smith, and A. G. Thakurta. Improved differential
280 privacy for SGD via optimal private linear operators on adaptive streams. In A. H. Oh, A. Agarwal,
281 D. Belgrave, and K. Cho, editors, *Advances in Neural Information Processing Systems*, 2022. URL
282 <https://openreview.net/forum?id=i9XrHJoyLqJ>.
- 283 D. Desfontaines. Differential privacy in (a bit) more detail. [https://desfontain.es/](https://desfontain.es/privacy/differential-privacy-in-more-detail.html)
284 [privacy/differential-privacy-in-more-detail.html](https://desfontain.es/privacy/differential-privacy-in-more-detail.html), 2018. Accessed 2023-
285 09-28.
- 286 Z. Ding, Y. Wang, G. Wang, D. Zhang, and D. Kifer. Detecting violations of differential privacy. In
287 *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security*,
288 pages 475–489, 2018.
- 289 C. Dwork and A. Roth. The algorithmic foundations of differential privacy. *Foundations and*
290 *Trends® in Theoretical Computer Science*, 9(3–4):211–407, 2014. ISSN 1551-305X. doi:
291 10.1561/0400000042. URL <http://dx.doi.org/10.1561/0400000042>.
- 292 C. Dwork, F. McSherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private data
293 analysis. In *Conference on Theory of Cryptography*, TCC ’06, pages 265–284, New York, NY,
294 USA, 2006.
- 295 A. C. Gilbert and A. McMillan. Property Testing for Differential Privacy. In *2018 56th Annual*
296 *Allerton Conference on Communication, Control, and Computing (Allerton)*, pages 249–258. IEEE,
297 2018.

- 298 N. Haim, G. Vardi, G. Yehudai, Michal Irani, and O. Shamir. Reconstructing training data from
299 trained neural networks. In A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho, editors, *Advances in*
300 *Neural Information Processing Systems*, 2022. URL [https://openreview.net/forum?](https://openreview.net/forum?id=Sxk8Bse3RKO)
301 [id=Sxk8Bse3RKO](https://openreview.net/forum?id=Sxk8Bse3RKO).
- 302 M. Jagielski, J. Ullman, and A. Oprea. Auditing differentially private machine learning: How private is
303 private sgd? In *Proceedings of the 34th International Conference on Neural Information Processing*
304 *Systems*, NIPS’20, Red Hook, NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.
- 305 M. Jagielski, S. Wu, A. Oprea, J. Ullman, and R. Geambasu. How to combine membership-
306 inference attacks on multiple updated machine learning models. *Proceedings on Privacy Enhancing*
307 *Technologies*, 3:211–232, 2023.
- 308 P. Kairouz, S. Oh, and P. Viswanath. The composition theorem for differential privacy. In F. Bach
309 and D. Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning*,
310 volume 37 of *Proceedings of Machine Learning Research*, pages 1376–1385, Lille, France, 07–09
311 Jul 2015. PMLR. URL <https://proceedings.mlr.press/v37/kairouz15.html>.
- 312 P. Kairouz, B. McMahan, S. Song, O. Thakkar, A. Thakurta, and Z. Xu. Practical and private (deep)
313 learning without sampling or shuffling. In *International Conference on Machine Learning*, pages
314 5213–5225. PMLR, 2021a.
- 315 P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles,
316 G. Cormode, R. Cummings, et al. Advances and open problems in federated learning. *Foundations*
317 *and Trends® in Machine Learning*, 14(1–2):1–210, 2021b.
- 318 X. Liu and S. Oh. Minimax Optimal Estimation of Approximate Differential Privacy on Neighboring
319 Databases. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- 320 F. Lu, J. Munoz, M. Fuchs, T. LeBlond, E. V. Zaresky-Williams, E. Raff, F. Ferraro, and B. Testa. A
321 general framework for auditing differentially private machine learning. In A. H. Oh, A. Agarwal,
322 D. Belgrave, and K. Cho, editors, *Advances in Neural Information Processing Systems*, 2022. URL
323 <https://openreview.net/forum?id=AKM3C3tsSx3>.
- 324 S. Maddock, A. Sablayrolles, and P. Stock. Canife: Crafting canaries for empirical privacy measure-
325 ment in federated learning, 2022. URL <https://arxiv.org/abs/2210.02912>.
- 326 B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y. Arcas. Communication-Efficient
327 Learning of Deep Networks from Decentralized Data. In A. Singh and J. Zhu, editors, *Proceed-*
328 *ings of the 20th International Conference on Artificial Intelligence and Statistics*, volume 54 of
329 *Proceedings of Machine Learning Research*, pages 1273–1282. PMLR, 20–22 Apr 2017. URL
330 <https://proceedings.mlr.press/v54/mcmahan17a.html>.
- 331 H. B. McMahan, D. Ramage, K. Talwar, and L. Zhang. Learning differentially private recurrent
332 language models. In *International Conference on Learning Representations*, 2018.
- 333 M. Nasr, S. Songi, A. Thakurta, N. Papemoti, and N. Carlin. Adversary instantiation: Lower bounds
334 for differentially private machine learning. In *2021 IEEE Symposium on Security and Privacy*
335 *(SP)*, pages 866–882. IEEE, 2021.
- 336 M. Nasr, J. Hayes, T. Steinke, B. Balle, F. Tramèr, M. Jagielski, N. Carlini, and A. Terzis. Tight
337 auditing of differentially private machine learning. *arXiv preprint arXiv:2302.07956*, 2023.
- 338 K. Pillutla, G. Andrew, P. Kairouz, H. B. McMahan, A. Oprea, and S. Oh. Unleashing the power of
339 randomization in auditing differentially private ml. *arXiv preprint arXiv:2305.18447*, 2023.
- 340 S. Reddi, Z. Charles, M. Zaheer, Z. Garrett, K. Rush, J. Konečný, S. Kumar, and H. B. McMahan.
341 Adaptive federated optimization. *arXiv preprint arXiv:2003.00295*, 2020.
- 342 R. Shokri, M. Stronati, C. Song, and V. Shmatikov. Membership inference attacks against machine
343 learning models. In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 3–18. IEEE, 2017.
- 344 T. Steinke. Composition of differential privacy privacy amplification by subsampling, 2022.

- 345 T. Steinke, M. Nasr, and M. Jagielski. Privacy auditing with one (1) training run. *arXiv preprint*
346 *arXiv:2305.08846*, 2023.
- 347 S. Yeom, I. Giacomelli, M. Fredrikson, and S. Jha. Privacy risk in machine learning: Analyzing the
348 connection to overfitting. In *2018 IEEE 31st Computer Security Foundations Symposium (CSF)*,
349 pages 268–282. IEEE, 2018.
- 350 S. Zanella-Beguelin, L. Wutschitz, S. Tople, A. Salem, V. Rühle, A. Paverd, M. Naseri, B. Köpf, and
351 D. Jones. Bayesian estimation of differential privacy. In *Proceedings of the 40th International*
352 *Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*,
353 pages 40624–40636. PMLR, 23–29 Jul 2023.
- 354 S. Zanella-Béguelin, L. Wutschitz, S. Tople, A. Salem, V. Rühle, A. Paverd, M. Naseri, B. Köpf,
355 and D. Jones. Bayesian estimation of differential privacy, 2022. URL [https://arxiv.org/](https://arxiv.org/abs/2206.05199)
356 [abs/2206.05199](https://arxiv.org/abs/2206.05199).

357 A Proofs of theorems from the main text

358 **Theorem 3.1.** *Let S be sampled uniformly from the unit sphere in \mathbb{R}^d , and let $\tau = \langle S, v \rangle / \|v\| \in$
359 $[-1, 1]$ be the cosine similarity between S and some arbitrary independent nonzero vector v . Then,
360 the probability density function of τ is*

$$f_d(\tau) = \frac{\Gamma(\frac{d}{2})}{\Gamma(\frac{d-1}{2})\sqrt{\pi}}(1 - \tau^2)^{\frac{d-3}{2}}.$$

361 *Proof.* Due to the rotational symmetry of the distribution of S , without loss of generality, we can
362 take v to be constant. First we describe the distribution of the angle $\theta \in [0, \pi]$ between s and v , then
363 change variables to get the distribution of its cosine τ . Consider the spherical cap of points on the
364 d -sphere with angle to v less or equal to θ , having $(d-1)$ -measure $A_d(\theta)$. The boundary of $A_d(\theta)$ is a
365 $(d-1)$ -sphere with radius $\sin \theta$ and $(d-2)$ -measure $M_d(\theta) = S_{d-1} \sin^{d-2} \theta$, where $S_d = 2\pi^{\frac{d}{2}}/\Gamma(\frac{d}{2})$
366 is the surface area of the unit d -sphere. (For example, the boundary of the 3-d spherical cap with
367 maximum angle θ is a circle (2-sphere) with radius $\sin \theta$ and circumference $2\pi \sin(\theta)$.) Normalizing
368 by the total area of the sphere S_d , the density of the angle is

$$\begin{aligned} \phi_d(\theta) &= S_d^{-1} \frac{d}{d\theta} A_d(\theta) \\ &= \left(\frac{2\pi^{\frac{d}{2}}}{\Gamma(\frac{d}{2})} \right)^{-1} \left(\frac{2\pi^{\frac{d-1}{2}}}{\Gamma(\frac{d-1}{2})} \right) \sin^{d-2} \theta \\ &= \frac{\Gamma(\frac{d}{2})}{\Gamma(\frac{d-1}{2})\sqrt{\pi}} \sin^{d-2} \theta. \end{aligned}$$

369 Now change variables to express it in terms of the angle cosine $\tau = \cos(\theta) \in [-1, 1]$:

$$\begin{aligned} f_d(\tau) &= \phi_d(\arccos \tau) \cdot \left| \frac{d}{d\tau} \arccos(\tau) \right| \\ &= \frac{\Gamma(\frac{d}{2})}{\Gamma(\frac{d-1}{2})\sqrt{\pi}} [\sin(\arccos \tau)]^{d-2} \left| -\frac{1}{\sqrt{1-\tau^2}} \right| \\ &= \frac{\Gamma(\frac{d}{2})}{\Gamma(\frac{d-1}{2})\sqrt{\pi}} \frac{(\sqrt{1-\tau^2})^{d-2}}{\sqrt{1-\tau^2}} \\ &= \frac{\Gamma(\frac{d}{2})}{\Gamma(\frac{d-1}{2})\sqrt{\pi}} (1 - \tau^2)^{\frac{d-3}{2}}. \end{aligned}$$

370 □

371 **Theorem 3.2.** *In the setting of Theorem 3.1, we have that $\tau\sqrt{d}$ converges in distribution to $\mathcal{N}(0, 1)$
372 as $d \rightarrow \infty$, i.e., $\lim_{d \rightarrow \infty} \mathbb{P}(\tau \leq \lambda/\sqrt{d}) = \mathbb{P}_{Z \sim \mathcal{N}(0,1)}(Z \leq \lambda)$.*

373 *Proof.* The distribution function of $t \in [-\sqrt{d}, \sqrt{d}]$ is

$$\begin{aligned}\hat{f}_d(t) &= \frac{1}{\sqrt{d}} \frac{\Gamma(\frac{d}{2})}{\Gamma(\frac{d-1}{2})\sqrt{\pi}} \left(1 - (t/\sqrt{d})^2\right)^{\frac{d-3}{2}} \\ &= \frac{\Gamma(\frac{d}{2})}{\Gamma(\frac{d-1}{2})\sqrt{\pi d}} \left(1 - t^2/d\right)^{\frac{d-3}{2}}.\end{aligned}$$

374 Taking the limit:

$$\begin{aligned}\lim_{d \rightarrow \infty} \hat{f}_d(t) &= \left(\lim_{d \rightarrow \infty} \frac{\Gamma(\frac{d}{2})}{\Gamma(\frac{d-1}{2})\sqrt{\pi d}}\right) \cdot \left(\lim_{d \rightarrow \infty} \left(1 - t^2/d\right)^{\frac{d}{2}}\right) \cdot \left(\lim_{d \rightarrow \infty} \left(1 - t^2/d\right)^{-\frac{3}{2}}\right) \\ &= \frac{1}{\sqrt{2\pi}} \cdot e^{-t^2/2} \cdot 1,\end{aligned}$$

375 where we have used the fact that $\frac{\Gamma(\frac{d}{2})}{\Gamma(\frac{d-1}{2})} \sim \sqrt{d/2}$.

376 □

377 **Lemma A.1.** *If τ is distributed according to the cosine angle distribution described in Theorem 3.1,*
378 *then $\text{Var}[\tau] = 1/d$.*

379 *Proof.* Let $x = (x_1, \dots, x_d)$ be uniform on the unit d -sphere. Then $x_1 = \langle x, e_1 \rangle$ has the required
380 distribution, where e_1 is the first standard basis vector. $\mathbb{E}[x_1]$ is zero, so we are interested in
381 $\text{Var}[x_1] = \mathbb{E}[x_1^2]$. Since $\sum_i x_i^2 = 1$, we have that $\mathbb{E}[\sum_i x_i^2] = \sum_i \mathbb{E}[x_i^2] = 1$. But all of the x_i
382 have the same distribution, so $\mathbb{E}[x_1^2] = 1/d$. □

383 **Theorem 3.3.** *For $d \in \mathbb{N}$, let $k = o(d)$, but $k = \omega(1)$. For $i = 1 \dots k$, let c_i sampled i.i.d.*
384 *from the unit sphere in d dimensions. Let $Z \sim \mathcal{N}(0; I_d)$. Let $\sigma > 0$, and define the mechanism*
385 *result $\rho = \sum_{j=1}^k c_j + \sigma Z$, and the cosine values $g_j = \frac{\langle c_j, \rho \rangle}{\|\rho\|}$. Write the empirical mean of the*
386 *cosines $\hat{\mu} = \frac{1}{k} \sum_{j=1}^k g_j$, and the empirical variance $\hat{\sigma}^2 = \frac{1}{k} \sum_{j=1}^k (g_j - \hat{\mu})^2$. Then as $d \rightarrow \infty$,*
387 *$\sqrt{d}\hat{\mu} \xrightarrow{p} 1/\sigma$ and $d\hat{\sigma}^2 \xrightarrow{p} 1$.*

388 *Proof.* Rewrite

$$\begin{aligned}\sqrt{d}\hat{\mu} &= \sqrt{d} \left(\frac{1}{k} \sum_i \frac{\langle c_i, \rho \rangle}{\|\rho\|} \right) \\ &= \left(\frac{\|\rho\|}{\sqrt{d}} \right)^{-1} \left(\frac{1}{k} \sum_i \langle c_i, \rho \rangle \right).\end{aligned}$$

389 We will show that $\frac{\|\rho\|^2}{d} \xrightarrow{p} \sigma^2$, while $\frac{1}{k} \sum_i \langle c_i, \rho \rangle \xrightarrow{p} 1$.

390 Note that $\|Z\|^2$ is Chi-squared distributed with mean d and variance $2d$. Also, for all $i \neq j$, $\langle c_i, c_j \rangle$
391 is distributed according to the cosine distribution discussed in Theorem 3.1 and Lemma A.1 with
392 mean 0 and variance $1/d$. Therefore,

$$\begin{aligned}\mathbb{E}\left[\frac{\|\rho\|^2}{d}\right] &= \frac{1}{d} \left\langle \sum_i c_i + \sigma Z, \sum_i c_i + \sigma Z \right\rangle \\ &= \frac{1}{d} \left(\sum_{i,j} \mathbb{E}[\langle c_i, c_j \rangle] + 2\sigma \sum_i \mathbb{E}[\langle Z, c_i \rangle] + \sigma^2 \mathbb{E}[\|Z\|^2] \right) \\ &= \frac{1}{d} (k + \sigma^2 d) \\ &= \sigma^2 + o(1).\end{aligned}$$

393 Next, note that all of the following dot products are pairwise uncorrelated: $\langle Z, Z \rangle$, $\langle Z, c_i \rangle$ (for all i),
 394 and $\langle c_i, c_j \rangle$ (for all i, j). Therefore the variance decomposes:

$$\begin{aligned} \text{Var} \left[\frac{\|\rho\|^2}{d} \right] &= \frac{1}{d^2} \left(\sum_{i,j} \text{Var}[\langle c_i, c_j \rangle] + 2\sigma^2 \sum_i \text{Var}[\langle Z, c_i \rangle] + \sigma^4 \text{Var}[\|Z\|^2] \right) \\ &= \frac{1}{d^2} \left(\frac{k(k-1)}{d} + 2k\sigma^2 + 2\sigma^4 d \right), \\ &= O(d^{-1}). \end{aligned}$$

395 Taken together, these imply that $\frac{\|\rho\|}{\sqrt{d}} \xrightarrow{p} \sigma$.

396 Now reusing many of the same calculations,

$$\mathbb{E} \left[\frac{1}{k} \sum_i \langle c_i, \rho \rangle \right] = \frac{1}{k} \left(k + \sum_{i \neq j} \mathbb{E}[\langle c_i, c_j \rangle] + \sigma \sum_i \mathbb{E}[\langle Z, c_i \rangle] \right) = 1,$$

397 and

$$\begin{aligned} \text{Var} \left[\frac{1}{k} \sum_i \langle c_i, \rho \rangle \right] &= \frac{1}{k^2} \left(\sum_{i \neq j} \text{Var}[\langle c_i, c_j \rangle] + \sigma^2 \sum_i \text{Var}[\langle Z, c_i \rangle] \right) \\ &= \frac{1}{k^2} \left(\frac{k(k-1)}{d} + k\sigma^2 \right) \\ &= o(1), \end{aligned}$$

398 which together imply that $\frac{1}{k} \sum_i \langle c_i, \rho \rangle \xrightarrow{p} 1$.

399 Now consider

$$\begin{aligned} d\hat{\sigma}^2 &= d \left(\frac{1}{k} \sum_i g_i^2 - \left(\frac{1}{k} \sum_i g_i \right)^2 \right) \\ &= \frac{d}{\|\rho\|^2} \left(\frac{1}{k} \sum_i \langle c_i, \rho \rangle^2 - \left(\frac{1}{k} \sum_i \langle c_i, \rho \rangle \right)^2 \right). \end{aligned}$$

400 We already have that $\frac{d}{\|\rho\|^2} \xrightarrow{p} \frac{1}{\sigma^2}$ and $\frac{1}{k} \sum_i \langle c_i, \rho \rangle \xrightarrow{p} 1$, so it will be sufficient to show
 401 $\frac{1}{k} \sum_i \langle c_i, \rho \rangle^2 \xrightarrow{p} 1 + \sigma^2$. Again using the uncorrelatedness of all pairs of dot products under
 402 consideration,

$$\begin{aligned} \mathbb{E} \left[\frac{1}{k} \sum_i \langle c_i, \rho \rangle^2 \right] &= \frac{1}{k} \sum_i \mathbb{E} \left[\left(\sum_j \langle c_i, c_j \rangle + \langle c_i, \sigma Z \rangle \right)^2 \right] \\ &= \frac{1}{k} \sum_i \left(\sum_{j,\ell} \mathbb{E}[\langle c_i, c_j \rangle \langle c_i, c_\ell \rangle] + 2 \sum_j \mathbb{E}[\langle c_i, c_j \rangle \langle c_i, \sigma Z \rangle] + \mathbb{E}[\langle c_i, \sigma Z \rangle^2] \right) \\ &= \frac{1}{k} \sum_i \left(\sum_{j,\ell} \mathbb{I}\{j = i, \ell = i\} + 0 + \sigma^2 \right) \\ &= 1 + \sigma^2. \end{aligned}$$

403 Also,

$$\text{Var} \left[\frac{1}{k} \sum_i \langle c_i, \rho \rangle^2 \right] = \frac{1}{k^2} \sum_i \left(\sum_{j,\ell} \text{Var}[\langle c_i, c_j \rangle \langle c_i, c_\ell \rangle] + 2 \sum_j \text{Var}[\langle c_i, c_j \rangle \langle c_i, \sigma Z \rangle] + \text{Var}[\langle c_i, \sigma Z \rangle^2] \right).$$

404 We'll bound each of these terms to show the sum is $o(1)$.

405 First, look at $\text{Var}[\langle c_i, c_j \rangle \langle c_i, c_\ell \rangle]$. If $i = j = \ell$, it is 0. If $j \neq \ell$,

$$\begin{aligned} \text{Var}[\langle c_i, c_j \rangle \langle c_i, c_\ell \rangle] &= \mathbb{E}[\langle c_i, c_j \rangle^2 \langle c_i, c_\ell \rangle^2] - \mathbb{E}[\langle c_i, c_j \rangle \langle c_i, c_\ell \rangle]^2 \\ &= \mathbb{E}[\mathbb{E}[\langle c_i, c_j \rangle^2 \langle c_i, c_\ell \rangle^2 \mid c_i]] \\ &= \mathbb{E}[\mathbb{E}[\langle c_i, c_j \rangle^2 \mid c_i] \mathbb{E}[\langle c_i, c_\ell \rangle^2 \mid c_i]] \\ &= \begin{cases} 1/d & j = i \text{ or } \ell = i, \\ 1/d^2 & j \neq i \text{ and } \ell \neq i, \end{cases} \end{aligned}$$

406 and if $j = \ell \neq i$,

$$\begin{aligned} \text{Var}[\langle c_i, c_j \rangle \langle c_i, c_\ell \rangle] &= \text{Var}[\langle c_i, c_j \rangle^2] \\ &= \mathbb{E}[\langle c_i, c_j \rangle^4] - \mathbb{E}[\langle c_i, c_j \rangle^2]^2 \\ &\leq \mathbb{E}[\langle c_i, c_j \rangle^4] \\ &\leq \mathbb{E}[\langle c_i, c_j \rangle^2] \\ &= 1/d. \end{aligned}$$

407 Together we have

$$\frac{1}{k^2} \sum_{i,j,\ell} \text{Var}[\langle c_i, c_j \rangle \langle c_i, c_\ell \rangle] \leq \frac{1}{k^2} \left(\frac{2k(k-1)}{d} + \frac{k(k-1)(k-2)}{d^2} + \frac{k(k-1)}{d} \right) = O(d^{-1}).$$

408 Now for $\text{Var}[\langle c_i, c_j \rangle \langle c_i, \sigma Z \rangle]$. If $i = j$, then it is σ^2 . If $i \neq j$,

$$\begin{aligned} \text{Var}[\langle c_i, c_j \rangle \langle c_i, \sigma Z \rangle] &= \mathbb{E}[\langle c_i, c_j \rangle^2 \langle c_i, \sigma Z \rangle^2] - \mathbb{E}[\langle c_i, c_j \rangle \langle c_i, \sigma Z \rangle]^2 \\ &= \mathbb{E}[\mathbb{E}[\langle c_i, c_j \rangle^2 \mid c_i] \mathbb{E}[\langle c_i, \sigma Z \rangle^2 \mid c_i]] \\ &= \sigma^2/d. \end{aligned}$$

409 So,

$$\frac{2}{k^2} \sum_{i,j} \text{Var}[\langle c_i, c_j \rangle \langle c_i, \sigma Z \rangle] = \frac{2}{k^2} \left(k\sigma^2 + \frac{k(k-1)\sigma^2}{d} \right) = o(1).$$

410 Finally, $\langle c_i, Z \rangle^2$ is Chi-squared distributed with one degree of freedom, so $\text{Var}[\langle c_i, \sigma Z \rangle^2] = 2\sigma^4$,
411 and

$$\frac{1}{k^2} \sum_i \text{Var}[\langle c_i, \sigma Z \rangle^2] = \frac{2\sigma^4}{k} = o(1).$$

412

□

413 B Algorithm for exact computation of ε comparing two Gaussian 414 distributions

415 In this section we give the details of the computation for estimating ε when $A(D)$ and $A(D')$ are
416 both Gaussian-distributed with different variances.

417 Let distribution under $A(D)$ be $P_1 = \mathcal{N}(\mu_1, \sigma_1^2)$ and the distribution under $A(D')$ be $P_2 =$
418 $\mathcal{N}(\mu_2, \sigma_2^2)$ with densities p_1 and p_2 respectively. Define $f_{P_1||P_2}(x) = \log \frac{p_1(x)}{p_2(x)}$. Now $Z_1 =$
419 $f_{P_1||P_2}(X_1)$ with $X_1 \sim P_1$ is the privacy loss random variable. Symmetrically define $Z_2 =$
420 $f_{P_2||P_1}(X_2)$ with $X_2 \sim P_2$.

421 From Steinke [2022] Prop. 7 we have that (ε, δ) -DP implies

$$\Pr[Z_1 > \varepsilon] - e^\varepsilon \Pr[-Z_2 > \varepsilon] \leq \delta \text{ and } \Pr[Z_2 > \varepsilon] - e^\varepsilon \Pr[-Z_1 > \varepsilon] \leq \delta.$$

422 Now we can compute

$$\begin{aligned}
f_{P_1||P_2}(x) &= \log \frac{p_1(x)}{p_2(x)} \\
&= \log \left(\frac{\sigma_2}{\sigma_1} \exp \left(-\frac{1}{2} \left[\frac{(x - \mu_1)^2}{\sigma_1^2} - \frac{(x - \mu_2)^2}{\sigma_2^2} \right] \right) \right) \\
&= \log \sigma_2 - \log \sigma_1 + \frac{(x - \mu_2)^2}{2\sigma_2^2} - \frac{(x - \mu_1)^2}{2\sigma_1^2} \\
&= ax^2 + bx + c
\end{aligned}$$

423 where

$$\begin{aligned}
a &= \frac{1}{2} \left(\frac{1}{\sigma_2^2} - \frac{1}{\sigma_1^2} \right), \\
b &= \frac{\mu_1}{\sigma_1^2} - \frac{\mu_2}{\sigma_2^2}, \\
\text{and } c &= \frac{1}{2} \left(\left(\frac{\mu_2}{\sigma_2} \right)^2 - \left(\frac{\mu_1}{\sigma_1} \right)^2 \right) + \log \sigma_2 - \log \sigma_1.
\end{aligned}$$

424 To compute $\Pr[Z_1 > \varepsilon]$, we need $\Pr[aX_1^2 + bX_1 + (c - \varepsilon) > 0]$ with $X_1 \sim P_1$. To do so, divide
425 the range of X_1 into intervals according to the zeros of $R(x) = ax^2 + bx + (c - \varepsilon)$. For example, if
426 R has roots $r_1 < r_2$ and a is positive, we can compute $\Pr[Z_1 > \varepsilon] = \Pr[X_1 < r_1] + \Pr[X_1 > r_2]$,
427 using the CDF of the Normal distribution. This requires considering a few cases, depending on the
428 sign of a and the sign of the determinant $b^2 - 4a(c - \varepsilon)$.

429 Now note that $f_{P_2||P_1} = -f_{P_1||P_2}$, so

$$\Pr[-Z_2 > \varepsilon] = \Pr[-f_{P_2||P_1}(X_2) > \varepsilon] = \Pr[aX_2^2 + bX_2 + (c - \varepsilon) > 0].$$

430 So the two events we are interested in ($Z_1 > \varepsilon$ and $-Z_2 > \varepsilon$) are the same, only when we compute
431 their probabilities according to P_1 vs. P_2 we use different values for μ and σ .

432 For numerical stability, the probabilities should be computed in the log domain. So we get

$$\begin{aligned}
\log \delta &\geq \log (\Pr[Z_1 > \varepsilon] - e^\varepsilon \Pr[-Z_2 > \varepsilon]) \\
&= \log \Pr[Z_1 > \varepsilon] + \log (1 - \exp(\varepsilon + \log \Pr[-Z_2 > \varepsilon] - \log \Pr[Z_1 > \varepsilon])).
\end{aligned}$$

433 Note it can happen that $\Pr[Z_1 > \varepsilon] < e^\varepsilon \Pr[-Z_2 > \varepsilon]$ in which case the corresponding bound is
434 invalid. A final trick we suggest for numerical stability is if $X \sim \mathcal{N}(\mu, \sigma^2)$ to use $\Pr(X < \mu; t, \sigma^2)$
435 in place of $\Pr(X > t; \mu, \sigma^2)$.

436 Now to determine ε at a given target δ , one can perform a line search over ε to find the value that
437 matches.

438 C Algorithms

439 Algorithms 1-3 illustrate the methods described in the main text.

440 D Supplementary experimental results on EMNIST dataset

441 In the main paper we presented results on the Stackoverflow federated word prediction task. Here
442 we present similar results on the EMNIST character recognition dataset. It contains 814k characters
443 written by 3383 users. The model is a CNN with 1.2M parameters. The users are shuffled and we
444 train for five epochs with 34 clients per round. The optimizers on client and server are both SGD,
445 with learning rates 0.031 and 1.0 respectively, and momentum of 0.9 on the server. The client batch
446 size is 16.

447 Table 3 shows the empirical epsilon estimates using either all model iterates or only the final model.
448 As with Stackoverflow next word prediction, using all iterates and a low amount of noise gives us

Algorithm 2 Privacy estimation via random canaries

1: Input: Client selection function clients , client training functions τ_i , canary selection function canaries , set of canary updates c_j , number of rounds T , initial parameters θ_0 , noise generator Z , ℓ_2 clip norm function S , privacy parameter δ , server learning rate η 2: for $t = 1, \dots, T$ do 3: $\rho = \vec{0}$ 4: for $i \in \mathbf{clients}(t)$ do 5: $\rho \leftarrow \rho + \text{CLIP}(\tau_i(\theta_{t-1}); S(t))$ 6: for $j \in \mathbf{canaries}(t)$ do 7: $\rho \leftarrow \rho + \text{PROJ}(c_j; S(t))$	8: $m = \mathbf{clients}(t) + \mathbf{canaries}(t) $ 9: $\theta_t \leftarrow \theta_{t-1} + \eta(\rho + Z(t))/m$ 10: for all canaries j do 11: $g_j \leftarrow \langle c_j, \theta_t \rangle / (\ c_j\ \cdot \ \theta_t\)$ 12: $\mu, \sigma \leftarrow \mathbf{mean}(\{g_j\}), \mathbf{std}(\{g_j\})$ 13: $\varepsilon \leftarrow \varepsilon(\mathcal{N}(0, 1/d) \parallel \mathcal{N}(\mu, \sigma^2); \delta)$ 14: function $\text{CLIP}(x; \kappa)$ 15: return $x \cdot \min(1, \kappa/\ x\)$ 16: function $\text{PROJ}(x; \kappa)$ 17: return $x \cdot \kappa/\ x\ $
---	---

Algorithm 3 Privacy estimation via random canaries using all iterates

1: Input: As in Algorithm 2, but with <i>unob-</i> <i>served</i> canary updates c_j^0 and <i>observed</i> canary updates c_j^1 . 2: for $t = 1, \dots, T$ do 3: $\rho = \vec{0}$ 4: for $i \in \mathbf{clients}(t)$ do 5: $\rho \leftarrow \rho + \text{CLIP}(\tau_i(\theta_{t-1}); S(t))$ 6: for $j \in \mathbf{canaries}(t)$ do 7: $\rho \leftarrow \rho + \text{PROJ}(c_j^1; S(t))$ 8: $m = \mathbf{clients}(t) + \mathbf{canaries}(t) $ 9: $\bar{\rho} \leftarrow (\rho + Z(t))/m$	10: for all canaries j do 11: $g_{t,j}^0 = \langle c_j^0, \bar{\rho} \rangle / (\ c_j^0\ \cdot \ \bar{\rho}\)$ 12: $g_{t,j}^1 = \langle c_j^1, \bar{\rho} \rangle / (\ c_j^1\ \cdot \ \bar{\rho}\)$ 13: $\theta_t \leftarrow \theta_{t-1} + \eta \bar{\rho}$ 14: for all canaries j do 15: $g_j^0 \leftarrow \max_t g_{t,j}^0$ 16: $g_j^1 \leftarrow \max_t g_{t,j}^1$ 17: $\mu_0, \sigma_0 \leftarrow \mathbf{mean}(\{g_j^0\}), \mathbf{std}(\{g_j^0\})$ 18: $\mu_1, \sigma_1 \leftarrow \mathbf{mean}(\{g_j^1\}), \mathbf{std}(\{g_j^1\})$ 19: $\varepsilon \leftarrow \varepsilon(\mathcal{N}(\mu_0, \sigma_0^2) \parallel \mathcal{N}(\mu_1, \sigma_1^2); \delta)$
--	---

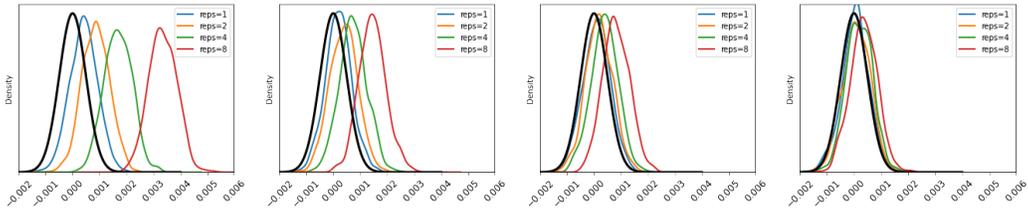


Figure 1: Density plots of cosine values with four values of noise corresponding to analytical epsilons (∞ , 300, 100, 30) and four values of canary repetitions (1, 2, 4, 8). The black curve in each plot is the pdf of the null distribution $\mathcal{N}(0, 1/d)$. With no noise ($\varepsilon = \infty$), the distributions are easily separable, with increasing separation for more canary repetitions. At higher levels of noise, distributions are less separable, even with several repetitions.

Noise	analytical ε	ε_{10} -all	ε_{est} -all	ε_{10} -final	ε_{est} -final
0.0	∞	6.25	48300	3.86	5.72
0.16	33.8	2.87	17.9	1.01	1.20
0.18	28.1	2.32	12.0	0.788	1.15
0.195	24.8	2.02	8.88	0.723	1.08
0.25	16.9	0.896	3.86	0.550	0.818
0.315	12.0	0.315	1.50	0.216	0.737

Table 3: Comparing ε estimates using all model deltas vs. using the final model only. ε_{10} is the empirical 95% lower bound from our modified Jagielski et al. [2020] method. The high values of ε_{est} -all indicate that membership inference is easy when the attacker has access to all iterates. On the other hand, when only the final model is observed, ε_{est} -final is far lower.

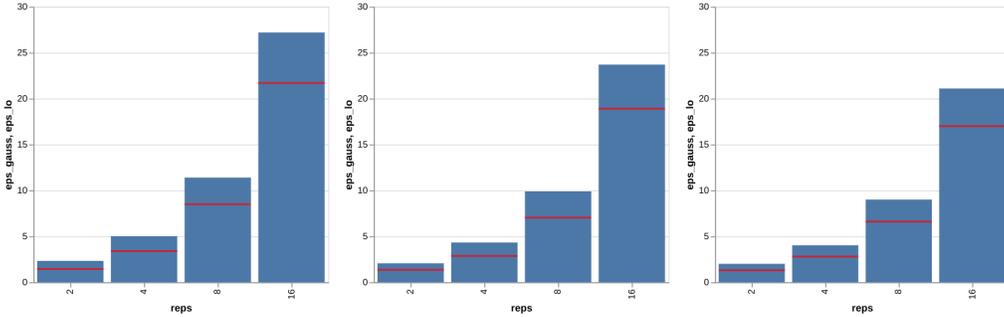


Figure 2: Blue bars are our ε_{est} and red ticks are the ε_{10} 95%-confidence lower bound for three noise multipliers (0.16, 0.18, 0.195) and four numbers of canary repetitions. Our estimate of epsilon increases sharply with the number of canary repetitions, confirming that limiting client participation improves privacy.

449 estimates close to the analytical upper bound, while using only the final model gives a much smaller
 450 estimate.

451 Figure 2 demonstrates the effect of increasing the number of canary repetitions for EMNIST. The
 452 results are qualitatively similar to the case of Stackoverflow.

453 E Experiments with multiple canary presentations

454 Here we highlight the ability of our method to estimate privacy when we vary not only the threat
 455 model, but also aspects of the training algorithm that may reasonably be expected to change privacy
 456 properties, but for which no tight analysis has been obtained. We consider presenting each canary a
 457 fixed multiple number of times, modeling the scenario in which clients are only allowed to check in
 458 for training every so often. In practice, a client need not participate in every period, but to obtain
 459 worst-case estimates, we present the canary in every period.

460 In Figure 1 we show kernel density estimation plots of the canary cosine sets. As the number of
 461 presentations increases in each plot, the distributions become more and more clearly separated.
 462 On the other hand as the amount of noise increases across the three plots, they converge to the
 463 null distribution. Also visible on this figure is that the distributions are roughly Gaussian-shaped,
 464 justifying the Gaussian approximation that is used in our estimation method. In Appendix G we
 465 give quantitative evidence for this observation. Finally we compare ε_{10} to our ε_{est} with multiple
 466 canary presentations in Figure 3. For each noise level, ε_{est} increases dramatically with increasing
 467 presentations, confirming our intuition that seeing examples multiple times dramatically reduces
 468 privacy.

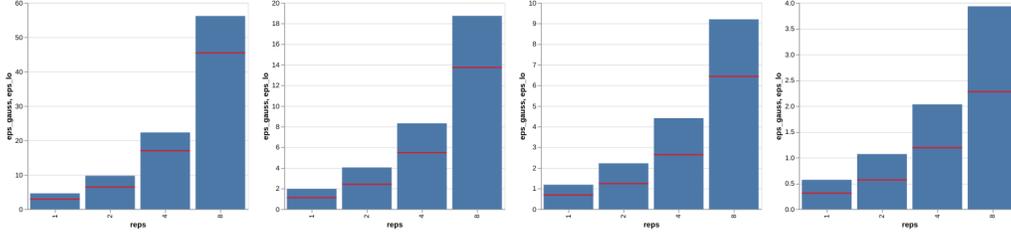


Figure 3: Blue bars are our ε_{est} and red ticks are the ε_{lo} 95%-confidence lower bound for four values of noise corresponding to analytical epsilons (∞ , 300, 100, 30) and four values of canary repetitions (1, 2, 4, 8). Note the difference of y-axis scales in each plot. Our estimate of epsilon increases sharply with the number of canary repetitions, confirming that limiting client participation improves privacy.

469 F Related Work

470 **Private federated learning.** DP Federated Averaging (DP-FedAvg) [McMahan et al., 2018] is a
 471 user-level DP version of the well-known Federated Averaging (FedAvg) algorithm [McMahan et al.,
 472 2017] for training ML models in a distributed fashion. In FedAvg, a central server interacts with
 473 a set of clients to train a global model iteratively over multiple rounds. In each round, the server
 474 sends the current global model to a subset of clients, who train local models using their training
 475 data, and send the model updates back to the server. The server aggregates the model updates via
 476 the Gaussian mechanism, in which each update is clipped to bound its ℓ_2 norm before averaging and
 477 adding Gaussian noise proportional to the clipping norm sufficient to mask the influence of individual
 478 users, and incorporates the aggregate update into the global model. DP-FedAvg can rely on privacy
 479 amplification from the sampling of clients at each round, but more sophisticated methods can handle
 480 arbitrary participation patterns [Kairouz et al., 2021a, Choquette-Choo et al., 2022].

481 **Privacy auditing.** Privacy auditing [Ding et al., 2018, Liu and Oh, 2019, Gilbert and McMillan,
 482 2018, Jagielski et al., 2020] provides techniques for empirically auditing the privacy leakage of
 483 an algorithm. The main technique used for privacy auditing is mounting a membership inference
 484 attack [Shokri et al., 2017, Yeom et al., 2018, Carlini et al., 2022], and translating the success of the
 485 adversary into an ε estimate using Eq. (1) directly.

486 Most privacy auditing techniques [Jagielski et al., 2020, Nasr et al., 2021, Lu et al., 2022, Zanella-
 487 Béguelin et al., 2022] have been designed for centralized settings, with the exception of CAN-
 488 IFE [Maddock et al., 2022], suitable for privacy auditing of federated learning deployments. CANIFE
 489 operates under a strong adversarial model, assuming knowledge of all intermediary model updates, as
 490 well as local model updates sent by a subset of clients in each round of training. CANIFE crafts data
 491 poisoning canaries adaptively, with the goal of generating model updates orthogonal to updates sent
 492 by other clients in each round. We argue that when the model dimensionality is sufficiently high, such
 493 crafting is unnecessary, since a randomly chosen canary update will already be essentially orthogonal
 494 to the true updates with high probability. CANIFE also computes a *per-round* privacy measure which
 495 it extrapolates into a measure for the entire training run by estimating an equivalent per-round noise
 496 $\hat{\sigma}_r$ and then composing the RDP of the repeated Poisson subsampled Gaussian mechanism. However
 497 in practice FL systems do not use Poisson subsampling due to the infeasibility of sampling clients i.i.d.
 498 at each round. Our method flexibly estimates the privacy loss in the context of arbitrary participation
 499 patterns, for example passing over the data in epochs, or the difficult-to-characterize *de facto* pattern
 500 of participation in a deployed system, which may include techniques intended to amplify privacy
 501 such as limits on client participation within temporal periods such as one day.

502 We empirically compare our approach with CANIFE in Appendix H and discuss the assumptions on
 503 the auditor’s knowledge and capabilities for all recent approaches (including ours) in Appendix I.

504 G Gaussianity of cosine statistics

505 To our knowledge, there is no way of confidently inferring that a set of samples comes from a given
 506 distribution, or even that they come from a distribution that is close to the given distribution in some
 507 metric. To quantify the error of our approximation of the cosine statistics with a Gaussian distribution,

Noise	1 rep	2 reps	3 reps	5 reps	10 reps
0	0.464	0.590	0.396	0.407	0.196
0.1023	0.976	0.422	0.340	0.432	0.116
0.2344	0.326	0.157	0.347	0.951	0.401

Table 4: Anderson statistics for each set of canary-cosine samples whose densities are shown in Figure 1. The Anderson test rejects at a 1% significance level if the statistic is greater than 1.088, and rejects at 15% significance if the statistic is greater than 0.574. Under the null hypothesis that all 15 (independent) distributions are Gaussian, we can compute that the probability of observing three or more values with Anderson test statistic greater than 0.574 is 68%.

508 we apply the Anderson test to each set of cosines in Table 4. It gives us some confidence to see
509 that a strong goodness-of-fit test cannot rule out that the distributions are Gaussian. This is a more
510 quantitative claim than visually comparing a histogram or empirical CDF, as is commonly done.

511 H Empirical comparison with CANIFE

512 As discussed in the main text, our method is significantly more general than the CANIFE method of
513 Maddock et al. [2022]. CANIFE periodically audits individual rounds to get a per-round $\hat{\epsilon}_r$, estimates
514 the noise for the round $\hat{\sigma}_r$ by inverting the computation of ϵ for the Gaussian mechanism, and uses
515 standard composition theorems to determine a final cumulative epsilon. Therefore if the assumptions
516 of those composition theorems do not strictly hold (for example, if clients are not sampled uniformly
517 and independently at each round), the estimate will be inaccurate. Also the step of crafting canaries
518 is model/dataset specific, and computationally expensive.

519 It is still interesting to see how the methods compare in the limited setting where CANIFE’s assump-
520 tions do hold. We trained a two-layer feedforward network on the fashion MNIST dataset. Following
521 experiments in Maddock et al. [2022], we used a canary design pool size of 512, took 2500 canary
522 optimization steps to find a canary example optimizing pixels and soft label, ran auditing every 100
523 rounds with 100 attack scores on each auditing round. We trained with a clip norm of 1.0 and noise
524 multiplier of 0.2 for one epoch with a batch size of 128, which corresponds to an analytical ϵ of 34.5.

525 CANIFE output a ϵ of mean 0.879 and standard deviation 0.124. Our method (using 1000 seen and
526 1000 unseen model canaries) estimated ϵ with a mean of 6.76 and std of 1.06, much closer to the
527 analytical bound.

528 I Comparison of assumptions and requirements of empirical privacy 529 estimation methods

530 As discussed in Section F, related work in privacy auditing/estimation rely on various assumptions on
531 the auditor’s knowledge and capability. Here we summarize the major differences.

		auditor controls	auditor receives
Central	Jagielski et al. [2020]	train data	final model
	Zanella-Beguelin et al. [2023]	train data	final model
	Pillutla et al. [2023]	train data	final model
	Steinke et al. [2023]	train data	final model
	Jagielski et al. [2023]	train data	intermediate models
	Nasr et al. [2023]	train data, privacy noise, minibatch	intermediate models
FL	Algorithm 2	client model update	final model
	Algorithm 3	client model update	intermediate models
	CANIFE [Maddock et al., 2022]	client sample, privacy noise, minibatch	intermediate models

Table 5: Examples of different assumptions in the literature. For each paper, we state the most relaxed condition the technique can be applied to, since they can be generalized in straightforward manner to scenarios with more strict assumptions on the auditor’s control and observation. Within each category of {central training, federated learning}, the lists are ordered from least to most strict assumptions.

Noise	runs/canaries	0.1	0.3	0.5	0.7	0.9
0.0496	1/1000	1.59	1.75	1.97	2.15	2.46
0.0496	10/100	1.65	1.80	1.92	2.16	2.32
0.0986	1/1000	0.81	1.06	1.18	1.33	1.54
0.0986	10/100	0.87	1.03	1.16	1.36	1.78

Table 6: Quantiles of $\hat{\epsilon}$ over fifty experiments using either one run with 1000 canaries or ten runs with 100 canaries each. For both noise multipliers, the distributions are very close.

532 Standard assumption in auditing centralized private training algorithm is a black-box setting where
533 the auditor only get to control the training data and observes the final model output. In practice, many
534 private training algorithms guarantee privacy under releasing all the intermediate model checkpoints.
535 One can hope to improve the estimate of privacy by using those check points as in [Jagielski et al.,
536 2023]. If the auditor can use information about how the minibatch sequence is drawn and the
537 distribution of the privacy noise, which is equivalent to assuming that the auditor controls the privacy
538 noise and the minibatch, one can further improve the estimates.

539 In the federated learning scenario, we assume canary client can return any model update. Note
540 that while CANIFE only controls the sample of the canary client and not the model update directly,
541 CANIFE utilises the Renyi-DP accountant with Poisson subsampling implemented via the Opacus
542 library, which is equivalent to the auditor fully controlling the sequence of minibatches (cohorts in
543 FL terminology). Further, the privacy noise is assumed to be independent spherical Gaussian, which
544 is equivalent to the auditor fully controlling the noise.

545 J Experiments comparing when multiple runs are used

546 In the limit of high model dimensionality, canaries are essentially mutually orthogonal, and therefore
547 they will interfere minimally with each other’s cosines to the model. In this section we give evidence
548 that even in the range of model dimensionalities explored in our experiments, including many canaries
549 in one run does not significantly perturb the estimated epsilon values. Ideally we would train 1000
550 models each with one canary to collect a set of truly independent statistics. However this is infeasible,
551 particularly if we want to perform the entire process multiple times to obtain confidence intervals.
552 Instead, we reduce the number of canaries per run by a factor of ten and train ten independent models
553 to collect a total of 1000 canary cosine statistics from which to estimate ϵ . We repeated the experiment
554 50 times for two different noise multipliers, which still amounts to training a total of 1000 models.
555 (Ten runs, two settings, fifty repetitions.)

556 The results on the stackoverflow dataset with the same setup as in Section 5 are shown in table 6.
557 We report the 0.1, 0.3, 0.5, 0.7 and 0.9 quantile of the distribution of $\hat{\epsilon}$ over 50 experiments. For
558 both noise multipliers, the distributions are quite close. Our epsilon estimates do not seem to vary
559 significantly even as the number of canaries per run varies by an order of magnitude.