
The Last Iterate Advantage: Empirical Auditing and Principled Heuristic Analysis of Differentially Private SGD

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We propose a simple heuristic privacy analysis of noisy clipped stochastic gradient
2 descent (DP-SGD) in the setting where only the last iterate is released and the in-
3 termediate iterates remain hidden. Namely, our heuristic assumes a linear structure
4 for the model.

5 We show experimentally that our heuristic is predictive of the outcome of privacy
6 auditing applied to various training procedures. Thus it can be used prior to training
7 as a rough estimate of the final privacy leakage. We also probe the limitations of
8 our heuristic by providing some artificial counterexamples where it underestimates
9 the privacy leakage.

10 The standard composition-based privacy analysis of DP-SGD effectively assumes
11 that the adversary has access to all intermediate iterates, which is often unrealistic.
12 However, this analysis remains the state of the art in practice. While our heuristic
13 does not replace a rigorous privacy analysis, it illustrates the large gap between
14 the best theoretical upper bounds and the privacy auditing lower bounds and sets a
15 target for further work to improve the theoretical privacy analyses.

16 1 Introduction

17 Differential privacy (DP) [DMNS06] defines a measure of how much private information from the
18 training data leaks through the output of an algorithm. The standard differentially private algorithm
19 for deep learning is DP-SGD [BST14; ACGMMTZ16], which differs from ordinary stochastic
20 gradient descent in two ways: the gradient of each example is clipped to bound its norm and then
21 Gaussian noise is added at each iteration.

22 The standard privacy analysis of DP-SGD is based on composition [BST14; ACGMMTZ16; Mir17;
23 Ste22; KJH20]. In particular, it applies to the setting where the privacy adversary has access to
24 all intermediate iterates of the training procedure. In this setting, the analysis is known to be tight
25 [NSTPC21; NHSBTJCT23]. However, in practice, potential adversaries rarely have access to the
26 intermediate iterates of the training procedure, rather they only have access to the final model. Access
27 to the final model can either be through queries to an API or via the raw model weights. The key
28 question motivating our work is the following.

29 Is it possible to obtain sharper privacy guarantees for DP-SGD when the adversary
30 only has access to the final model, rather than all intermediate iterates?

31 1.1 Background & Related Work

32 The question above has been studied from two angles: Theoretical upper bounds, and privacy auditing
33 lower bounds. Our goal is to shed light on this question from a third angle via principled heuristics.

34 A handful of theoretical analyses [FMTT18; CYS21; YS22; AT22; BSA24] have shown that asymptotically the privacy guarantee of the last iterate of DP-SGD can be far better than the standard
35 composition-based analysis that applies to releasing all iterates. In particular, as the number of
36 iterations increases, these analyses give a privacy guarantee that converges to a constant (depending
37 on the loss function and the scale of the noise), whereas the standard composition-based analysis
38 would give a privacy guarantee that increases forever. Unfortunately, these theoretical analyses
39 are only applicable under strong assumptions on the loss function, such as (strong) convexity and
40 smoothness. We lack an understanding of how well they reflect the “real” privacy leakage.

42 Privacy auditing [JUO20; DWWZK18; BGDCTV18; SNJ23; TTSSJC22; ZBWTSRPNK22] complements
43 theoretical analysis by giving empirical lower bounds on the privacy leakage. Privacy
44 auditing works by performing a membership inference attack [SSSS17; HSRDTMPSNC08; SOJH09;
45 DSSUV15]. That is, it constructs neighbouring inputs and demonstrates that the corresponding
46 output distributions can be distinguished well enough to imply a lower bound on the differential
47 privacy parameters. In practice, the theoretical privacy analysis may give uncomfortably large values
48 for the privacy leakage (e.g., $\epsilon > 10$); in this case, privacy auditing may be used as evidence that
49 the “real” privacy leakage is lower. There are settings where the theoretical analysis is matched by
50 auditing, such as when all intermediate results are released [NSTPC21; NHSBTJCT23]. However,
51 despite significant work on privacy auditing and membership inference [CCNSTT22; BTRKMW24;
52 WBKGGGG23; LF20; SDSOJ19; ZLS23], a large gap remains between the theoretical upper bounds
53 and the auditing lower bounds [AKOOMS23; NHSBTJCT23] when only the final parameters are
54 released. This observed gap is the starting point for our work.

55 1.2 Our Contributions

56 We propose a *heuristic* privacy analysis of DP-SGD in the setting where only the final iterate is
57 released. Our experiments demonstrate that this heuristic analysis consistently provides an upper
58 bound on the privacy leakage measured by privacy auditing tools in realistic deep learning settings.

59 Our heuristic analysis corresponds to a worst-case theoretical analysis under the assumption that the
60 loss functions are linear. This case is simple enough to allow for an exact privacy analysis whose
61 parameters can be computed numerically (Theorem 1). Our consideration of linear losses is built
62 on the observation that current auditing techniques achieve the highest ϵ values when the gradients
63 of the canaries – that is, the examples that are included or excluded to test the privacy leakage – are
64 fixed and independent from the gradients of the other examples. This is definitely the case for linear
65 losses; the linear assumption thus allows us to capture the setting where current attacks are most
66 effective. Linear loss functions are also known to be the worst case for the non-sampled (i.e., full
67 batch) case; see Appendix B. Assuming linearity is unnatural from an optimization perspective, as
68 there is no minimizer. But, from a privacy perspective, we show that it captures the state of the art.

69 We also probe the limitations of our heuristic and give some artificial counterexamples where it
70 underestimates empirical privacy leakage. One class of counterexamples exploits the presence of a
71 regularizer. Roughly, the regularizer partially zeros out the noise that is added for privacy. However,
72 the regularizer also partially zeros out the signal of the canary gradient. These two effects are almost
73 balanced, which makes the counterexample very delicate. In a second class of counterexamples, the
74 data is carefully engineered so that the final iterate effectively encodes the entire trajectory, in which
75 case there is no difference between releasing the last iterate and all iterates.

76 **Implications:** Heuristics cannot replace rigorous theoretical analyses. However, our heuristic can
77 serve as a target for future improvements to both privacy auditing as well as theoretical analysis. For
78 privacy auditing, matching or exceeding our heuristic is a more reachable goal than matching the
79 theoretical upper bounds, although our experimental results show that even this would require new
80 attacks. When theoretical analyses fail to match our heuristic, we should identify why there is a gap,
81 which builds intuition and could point towards further improvements.

82 Given that privacy auditing is computationally intensive and difficult to perform correctly [AZT24],
83 we believe that our heuristic can also be valuable in practice. In particular, our heuristic can be used
84 prior to training (e.g., during hyperparameter selection) to predict the outcome of privacy auditing
85 when applied to the final model. (This is a similar use case to scaling laws.)

86 2 Linearized Heuristic Privacy Analysis

87 Theorem 1 presents our heuristic differential
88 privacy analysis of DP-SGD (which we
89 present in Algorithm 1 for completeness;
90 note that we include a regularizer r whose
91 gradient is *not* clipped, because it does not
92 depend on the private data \mathbf{x}). We con-
93 sider Poisson subsampled minibatches and
94 add/remove neighbours, as is standard in
95 the differential privacy literature.

96 Our analysis takes the form of a conditional
97 privacy guarantee. Namely, under the as-
98 sumption that the loss and regularizer are
99 linear, we obtain a fully rigorous differen-
100 tial privacy guarantee. The heuristic is to
101 apply this guarantee to loss functions that
102 are not linear (such as those that arise in
103 deep learning applications). Our thesis is
104 that, in most cases, the conclusion of the
105 theorem is still a good approximation, even
106 when the assumption does not hold.

107 Recall that a function $\ell : \mathbb{R}^d \rightarrow \mathbb{R}$ is linear
108 if there exist $\alpha \in \mathbb{R}^d$ and $\beta \in \mathbb{R}$ such that
109 $\ell(\mathbf{m}) = \langle \alpha, \mathbf{m} \rangle + \beta$ for all \mathbf{m} .

110 **Theorem 1** (Privacy of DP-SGD for linear
111 losses). *Let $\mathbf{x}, T, q, \eta, \sigma, \ell, r$ be as in Algo-*
112 *rithm 1. Assume r and $\ell(\cdot, x)$, for every*
113 *$x \in \mathcal{X}$, are linear.*

114 *Letting*

$$P := \text{Binomial}(T, q) + \mathcal{N}(0, \sigma^2 T), Q := \mathcal{N}(0, \sigma^2 T), \quad (1)$$

115 *DP-SGD with `last_iterate_only` satisfies (ε, δ) -differential privacy with $\varepsilon \geq 0$ arbitrary and*

$$\delta = \delta_{T,q,\sigma}(\varepsilon) := \max\{H_{e^\varepsilon}(P, Q), H_{e^\varepsilon}(Q, P)\}. \quad (2)$$

116 *Here, H_{e^ε} denotes the e^ε -hockey-stick-divergence $H_{e^\varepsilon}(P, Q) := \sup_S P(S) - e^\varepsilon Q(S)$.*

117 Equation 1 gives us a value of the privacy failure probability parameter δ . But it is more natural to
118 work with the privacy loss bound parameter ε , which can be computed by inverting the formula:

$$\varepsilon_{T,q,\sigma}(\delta) := \min\{\varepsilon \geq 0 : \delta_{T,q,\sigma}(\varepsilon) \leq \delta\}. \quad (3)$$

119 Both $\delta_{T,q,\sigma}(\varepsilon)$ and $\varepsilon_{T,q,\sigma}(\delta)$ can be computed using existing open-source DP accounting libraries
120 [Goo20]. We also provide a self-contained & efficient method for computing them in Appendix A.

121 The proof of Theorem 1 is deferred to Appendix A, but we sketch the main ideas: Under the linearity
122 assumption, the output of DP-SGD is just a sum of the gradients and noises. We can reduce to
123 dimension $d = 1$, since the only relevant direction is that of the gradient of the canary¹ (which is
124 constant). We can also ignore the gradients of the other examples. Thus, by rescaling, the worst case
125 pair of output distributions can be represented as in Equation 1. Namely, $Q = \sum_{t=1}^T \xi_t$ is simply the
126 noise $\xi_t \leftarrow \mathcal{N}(0, \sigma^2)$ summed over T iterations; this corresponds to the case where the canary is
127 excluded. When the canary is included, it is sampled with probability q in each iteration and thus
128 the total number of times it is sampled over T iterations is $\text{Binomial}(T, q)$. Thus P is the sum of the
129 contributions of the canary and the noise. Finally the definition of differential privacy lets us compute
130 ε and δ from this pair of distributions. Tightness follows from the fact that there exists a loss function
131 and pair of inputs such that the corresponding outputs of DP-SGD matches the pair P and Q .

¹The *canary* refers to the individual datapoint that is added or removed between neighbouring datasets. This terminology is used in the privacy auditing/attacks literature inspired on the expression ‘‘canary in a coalmine.’’

Algorithm 1 Noisy Clipped Stochastic Gradient Descent (DP-SGD) [BST14; ACGMMTZ16]

function DP-SGD($\mathbf{x} \in \mathcal{X}^n, T \in \mathbb{N}, q \in [0, 1], \eta \in (0, \infty), \sigma \in (0, \infty), \ell : \mathbb{R}^d \times \mathcal{X} \rightarrow \mathbb{R}, r : \mathbb{R}^d \rightarrow \mathbb{R}$)

Initialize model $\mathbf{m}_0 \in \mathbb{R}^d$.

for $t = 1 \dots T$ **do**

Sample minibatch $B_t \subseteq [n]$ including each element independently with probability q .

Compute gradients of the loss $\nabla_{\mathbf{m}_{t-1}} \ell(\mathbf{m}_{t-1}, x_i)$ for all $i \in B_t$ and of the regularizer $\nabla_{\mathbf{m}_{t-1}} r(\mathbf{m}_{t-1})$.

Clip loss gradients: $\text{clip}(\nabla_{\mathbf{m}_{t-1}} \ell(\mathbf{m}_{t-1}, x_i)) := \frac{\nabla_{\mathbf{m}_{t-1}} \ell(\mathbf{m}_{t-1}, x_i)}{\max\{1, \|\nabla_{\mathbf{m}_{t-1}} \ell(\mathbf{m}_{t-1}, x_i)\|_2\}}$.

Sample noise $\xi_t \leftarrow \mathcal{N}(0, \sigma^2 I_d)$.

Update

$$\mathbf{m}_t = \mathbf{m}_{t-1} - \eta \cdot \left(\frac{\sum_{i \in B_t} \text{clip}(\nabla_{\mathbf{m}_{t-1}} \ell(\mathbf{m}_{t-1}, x_i))}{+\nabla_{\mathbf{m}_{t-1}} r(\mathbf{m}_{t-1})} + \xi_t \right).$$

end for

if `last_iterate_only` **then**

return \mathbf{m}_T

else if `intermediate_iterates` **then**

return $\mathbf{m}_0, \mathbf{m}_1, \dots, \mathbf{m}_{T-1}, \mathbf{m}_T$

end if

end function

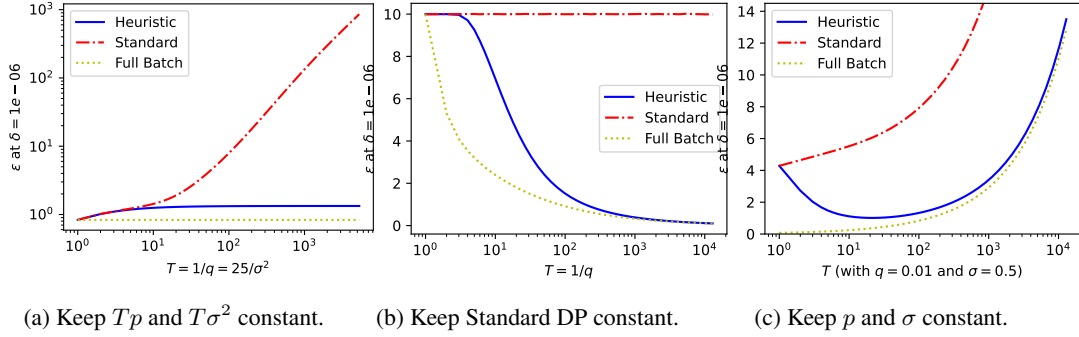


Figure 1: Comparison of our heuristic to baselines in various parameter regimes. Horizontal axis is number of iterations T and vertical axis is ϵ such that we have $(\epsilon, 10^{-6})$ -DP.

132 2.1 Baselines

133 In addition to privacy auditing, we compare our heuristic to two different baselines in Figure 1.
 134 The first is the standard, composition-based analysis. We use the open-source library from Google
 135 [Goo20], which computes a tight DP guarantee for DP-SGD with `intermediate_iterates`. Be-
 136 cause DP-SGD with `intermediate_iterates` gives the adversary more information than with
 137 `last_iterate_only`, this will always give at least as large an estimate for ϵ as our heuristic.

138 We also consider approximating DP-SGD by full batch DP-GD. That is, set $q = 1$ and rescale the
 139 learning rate η and noise multiplier σ to keep the expected step and privacy noise variance constant:

$$\underbrace{\text{DP-SGD}(\mathbf{x}, T, q, \eta, \sigma, \ell, r)}_{\text{batch size } \approx nq, T \text{ iterations, } Tq \text{ epochs}} \approx \underbrace{\text{DP-SGD}(\mathbf{x}, T, 1, \eta \cdot q, \sigma/q, \ell, r)}_{\text{batch size } n, T \text{ iterations, } T \text{ epochs}}. \quad (4)$$

140 The latter algorithm is full batch DP-GD since at each step it includes each data point in the batch
 141 with probability 1. Since full batch DP-GD does not rely on privacy amplification by subsampling,
 142 it is much easier to analyze its privacy guarantees. Interestingly, there is no difference between
 143 full batch DP-GD with `last_iterate_only` and with `intermediate_iterates`; see Appendix B.
 144 Full batch DP-GD generally has better privacy guarantees than the corresponding minibatch DP-SGD
 145 and so this baseline usually (but not always) gives smaller values for the privacy leakage ϵ than our
 146 heuristic. In practice, full batch DP-GD is too computationally expensive to run. But we can use it as
 147 an idealized comparison point for the privacy analysis.

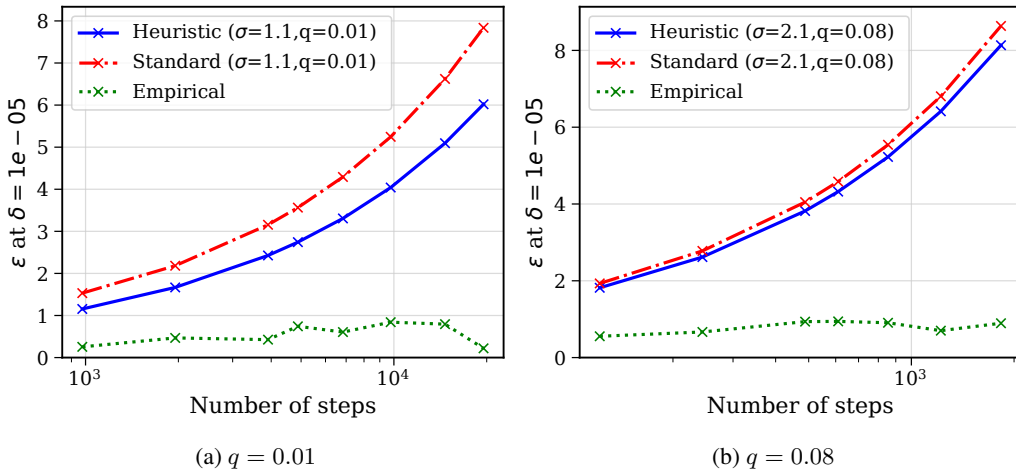


Figure 2: Black-box gradient space attacks fail to achieve tight auditing when other data points are sampled from the data distribution. Heuristic and standard bounds diverge from empirical results, indicating the attack’s ineffectiveness. This contrasts with previous work which tightly auditing with access to intermediate updates.

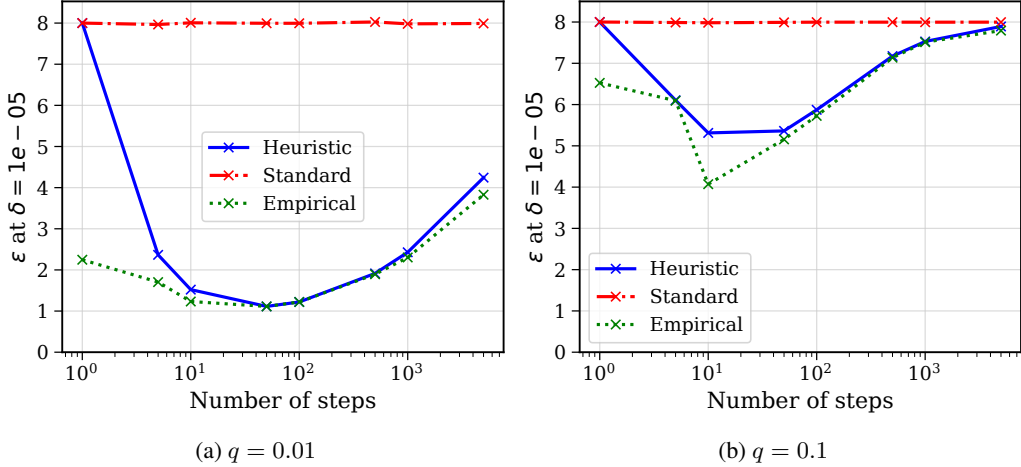


Figure 3: For gradient space attacks with adversarial datasets, the empirical epsilon (ε) closely tracks the final epsilon except for at small step counts, where distinguishing is more challenging. This is evident at both subsampling probability values we study ($q = 0.01$ and $q = 0.1$).

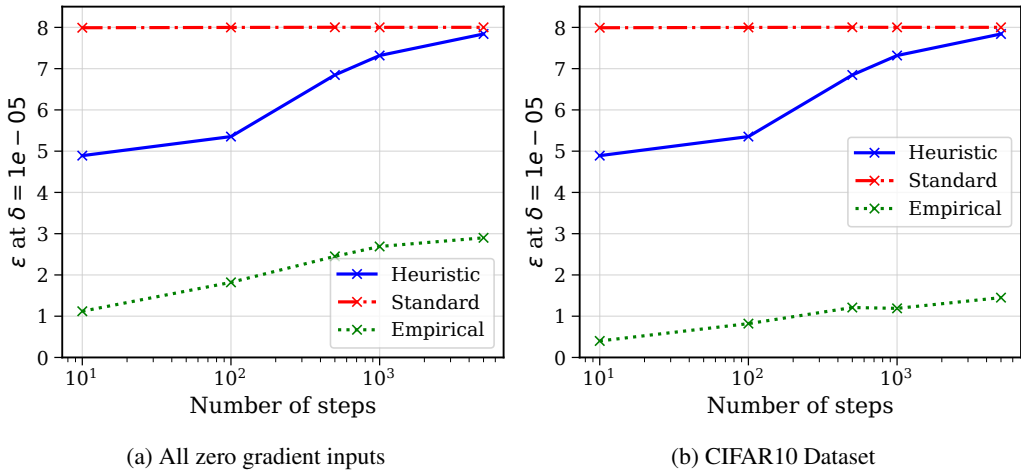


Figure 4: Input space attacks show promising results with both natural and blank image settings, although blank images have higher attack success. These input space attacks achieve tighter results than gradient space attacks in the natural data setting, in contrast to findings from prior work.

148 3 Empirical Evaluation via Privacy Auditing

149 **Setup:** We follow the construction of Nasr, Song, Thakurta, Papernot, and Carlini [NSTPC21] where
 150 we have 3 entities, adversarial crafter, model trainer, and distinguisher. In this paper, we assume
 151 the distinguisher only has access the final iteration of the model parameters. We use the CIFAR10
 152 dataset [Ale09] with a WideResNet model [ZK16] unless otherwise specified; in particular, we follow
 153 the training setup of De, Berrada, Hayes, Smith, and Balle [DBHSB22], where we train and audit
 154 a model with 79% test accuracy and, using the standard analysis, ($\varepsilon = 8, \delta = 10^{-5}$)-DP. For each
 155 experiment we trained 512 CIFAR10 models with and without the canary (1024 total). To compute
 156 the empirical lower bounds we use the PLD approach with Clopper-Pearson confidence intervals
 157 used by Nasr, Hayes, Steinke, Balle, Tramèr, Jagielski, Carlini, and Terzis [NHSBTJCT23]. Here we
 158 assume the adversary knows the sampling rate and the number of iterations and is only estimating the
 159 noise multiplier used in DP-SGD, from which the reported privacy parameters (ε and δ) are derived.

160 3.1 Experimental Results

161 We implement state-of-the-art attacks from prior work [NSTPC21; NHSBTJCT23]. These attacks
162 heavily rely on the intermediate steps and, as a result, do not achieve tight results. In the next
163 section, we design specific attacks for our heuristic privacy analysis approach to further understand
164 its limitations and potential vulnerabilities. We used Google Cloud A2-megagpu-16g machines with
165 16 Nvidia A100 40GB GPUs. Overall, we use roughly 33,000 GPU hours for our experiments.

166 **Gradient Space Attack:** The most powerful attacks in prior work are gradient space attacks where
167 the adversary injects a malicious gradient directly into the training process, rather than an example;
168 prior work has shown that this attack can produce tight lower bounds, independent of the dataset
169 and model used for training [NHSBTJCT23]. However, these previous attacks require access to all
170 intermediate training steps to achieve tight results. Here, we use canary gradients in two settings: one
171 where the other data points are non-adversarial and sampled from the real training data, and another
172 where the other data points are designed to have very small gradients (≈ 0). This last setting was
173 shown by [NSTPC21] to result in tighter auditing. In all attacks, we assume the distinguisher has
174 access to all adversarial gradient vectors. For malicious gradients, we use Dirac gradient canaries,
175 where gradient vectors consist of zeros in all but a single index. In both cases, the distinguishing test
176 measures the dot product of the final model checkpoint and the gradient canary.

177 Figure 2 summarizes the results for the non-adversarial data setting, with other examples sampled
178 from the true training data. In this experiment, we fix noise magnitude and subsampling probability,
179 and run for various numbers of training steps. While prior work has shown tight auditing in this
180 setting, we find an adversary without access to intermediate updates obtains much weaker attacks.
181 Indeed, auditing with this strong attack results even in much lower values than the heuristic outputs.

182 Our other setting assumes the other data points are maliciously chosen. We construct an adversarial
183 “dataset” of $m + 1$ gradients, m of which are zero, and one gradient is constant (with norm equal to
184 the clipping norm), applying gradients directly rather than using any examples. As this experiment
185 does not require computing gradients, it is very cheap to run more trials, so we run this procedure
186 $N = 100,000$ times with the gradient canary, and N times without it, and compute an empirical
187 estimate for ε with these values. We plot the results of this experiment in Figure 3 together with
188 the ε output by the theoretical analysis and the heuristic, fixing the subsampling probability and
189 varying the number of update steps. We adjust the noise parameter to ensure the standard theoretical
190 analysis produces a fixed ε bound. The empirical measured ε is close to the heuristic ε except for
191 when training with very small step counts: we expect this looseness to be the result of statistical
192 effects, as lower step counts have higher relative variance at a fixed number of trials.

193 **Input Space Attack:** In practice, adversaries typically cannot insert malicious gradients freely in
194 training steps. Therefore, we also study cases where the adversary is limited to inserting malicious
195 inputs into the training set. Label flip attacks are one of the most successful approaches used to audit
196 DP machine learning models in prior work [NHSBTJCT23; SNJ23]. For input space attacks, we use
197 the loss of the malicious input as a distinguisher. Similar to our gradient space attacks, we consider
198 two settings for input space attacks: one where other data points are correctly sampled from the
199 dataset, and another where the other data points are blank images.

200 Figure 4 summarizes the results for this setting. Comparing to Figure 2, input space attacks achieve
201 tighter results than gradient space attacks. This finding is in stark contrast to prior work. The reason
202 is that input space attacks do not rely on intermediate iterates, so they transfer well to our setting.

203 In all the cases discussed so far, the empirical results for both gradient and input attacks fall below
204 the heuristic analysis and do not violate the upper bounds based on the underlying assumptions. This
205 suggests that the heuristic might serve as a good indicator for assessing potential vulnerabilities.
206 However, in the next section, we delve into specific attack scenarios that exploit the assumptions used
207 in the heuristic analysis to create edge cases where the heuristic bounds are indeed violated.

208 4 Counterexamples

209 We now test the limits of our heuristic by constructing some artificial counterexamples. That is, we
210 construct inputs to DP-SGD with `last_iterate_only` such that the true privacy loss exceeds the
211 bound given by our heuristic. While we do not expect the contrived structures of these examples to

212 manifest in realistic learning settings, they highlight the difficulties of formalizing settings where the
 213 heuristic gives a provable upper bound on the privacy loss.

214 4.1 Warmup: Zeroing Out The Model Weights

215 We begin by noting the counterintuitive fact that our heuristic $\varepsilon_{T,q,\sigma}(\delta)$ is *not* always monotone in
 216 the number of steps T when the other parameters σ, q, δ are kept constant. This is shown in Figure 1c.
 217 More steps means there is both more noise and more signal from the gradients; these effects partially
 218 cancel out, but the net effect can be non-monotone.

219 We can use a regularizer $r(\mathbf{m}) = \|\mathbf{m}\|_2^2/2\eta$ so that $\eta \cdot \nabla_{\mathbf{m}} r(\mathbf{m}) = \mathbf{m}$. This regularizer zeros out the
 220 model from the previous step, i.e., the update of DP-SGD becomes

$$\mathbf{m}_t = \mathbf{m}_{t-1} - \eta \cdot \left(\sum_{i \in B_t} \text{clip}(\nabla_{\mathbf{m}_{t-1}} \ell(\mathbf{m}_{t-1}, x_i)) + \nabla_{\mathbf{m}_{t-1}} r(\mathbf{m}_{t-1}) + \xi_t \right) \quad (5)$$

$$= \eta \cdot \sum_{i \in B_t} \text{clip}(\nabla_{\mathbf{m}_{t-1}} \ell(\mathbf{m}_{t-1}, x_i)) + \xi_t. \quad (6)$$

221 This means that the last iterate \mathbf{m}_T is effectively the result of only a single iteration of DP-SGD. In
 222 particular, it will have a privacy guarantee corresponding to one iteration. Combining this regularizer
 223 with a linear loss and a setting of the parameters T, q, σ, δ such that the privacy loss is non-monotone
 224 – i.e., $\varepsilon_{T,q,\sigma}(\delta) < \varepsilon_{1,q,\sigma}(\delta)$ – yields a counterexample.

225 In light of this counterexample, in the next subsection, we benchmark our counterexample against
 226 sweeping over smaller values of T . I.e., we consider $\max_{t \leq T} \varepsilon_{t,q,\sigma}(\delta)$ instead of simply $\varepsilon_{T,q,\sigma}(\delta)$.

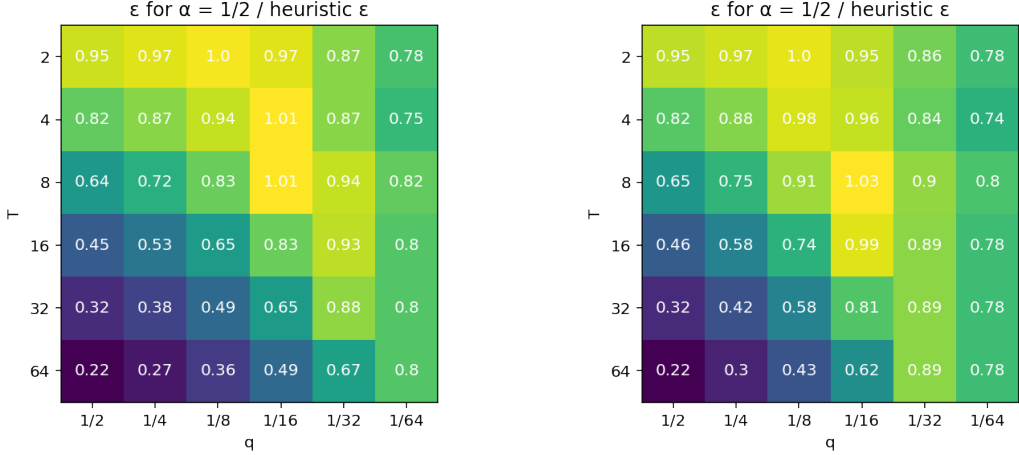
227 4.2 Linear Loss + Quadratic Regularizer

228 Consider running DP-SGD in one dimension (i.e., $d = 1$) with a linear loss $\ell(\mathbf{m}, x) = \mathbf{m}x$ for
 229 the canary and a quadratic regularizer $r(\mathbf{m}) = \frac{1}{2}\alpha\mathbf{m}^2$, where $\alpha \in [0, 1]$ and $x \in [-1, 1]$ and we
 230 use learning rate $\eta = 1$. With sampling probability q , after T iterations the privacy guarantee
 231 is equivalent to distinguishing $Q := \mathcal{N}(0, \hat{\sigma}^2)$ and $P := \mathcal{N}(\sum_{i \in [T]} (1 - \alpha)^{i-1} \text{Bernoulli}(q), \hat{\sigma}^2)$,
 232 where $\hat{\sigma}^2 := \sigma^2 \sum_{i \in [T]} (1 - \alpha)^{2(i-1)}$. When $\alpha = 0$, this retrieves linear losses. When $\alpha = 1$, this
 233 corresponds to distinguishing $\mathcal{N}(0, \hat{\sigma}^2)$ and $\mathcal{N}(\text{Bernoulli}(q), \hat{\sigma}^2)$ or, equivalently, to distinguishing
 234 linear losses after $T = 1$ iteration. If we maximize our heuristic over the number of iterations $\leq T$,
 235 then our heuristic is tight for the extremes $\alpha \in \{0, 1\}$.

236 A natural question is whether the worst-case privacy guarantee on this quadratic is always given by
 237 $\alpha \in \{0, 1\}$. Perhaps surprisingly, the answer is no: we found that for $T = 3, q = 0.1, \sigma = 1, \alpha = 0$,
 238 DP-SGD is $(2.222, 10^{-6})$ -DP. For $\alpha = 1$ instead DP-SGD is $(2.182, 10^{-6})$ -DP. However, for
 239 $\alpha = 0.5$ instead the quadratic loss does not satisfy $(\varepsilon, 10^{-6})$ -DP for $\varepsilon < 2.274$.

240 However, this violation is small, which suggests our heuristic is still a reasonable for this class of
 241 examples. To validate this, we consider a set of values for the tuple (T, q, σ) . For each setting
 242 of T, q, σ , we compute $\max_{t \leq T} \varepsilon_{t,q,\sigma}(\delta)$ at $\delta = 10^{-6}$. We then compute ε for the linear loss
 243 with quadratic regularizer example with $\alpha = 1/2$ in the same setting. Since the support of the
 244 random variable $\sum_{i \in [T]} (1 - \alpha)^{i-1} \text{Bernoulli}(q)$ has size 2^T for $\alpha = 1/2$, computing exact ε for
 245 even moderate T is computationally intensive. Instead, let X be the random variable equal to
 246 $\sum_{i \in [T]} (1 - \alpha)^{i-1} \text{Bernoulli}(q)$, except we round up values in the support which are less than .0005
 247 up to .0005, and then round each value in the support up to the nearest integer power of 1.05. We then
 248 compute an exact ε for distinguishing $\mathcal{N}(0, \hat{\sigma}^2)$ vs $\mathcal{N}(X, \hat{\sigma}^2)$. By Lemma 4.5 of Choquette-Choo,
 249 Ganesh, Steinke, and Thakurta [CCGST24], we know that distinguishing $\mathcal{N}(0, \hat{\sigma}^2)$ vs. $\mathcal{N}(\sum_{i \in [T]} (1 -$
 250 $\alpha)^{i-1} \text{Bernoulli}(q), \hat{\sigma}^2)$ is no harder than distinguishing $\mathcal{N}(0, \hat{\sigma}^2)$ vs $\mathcal{N}(X, \hat{\sigma}^2)$, and since we increase
 251 the values in the support by no more than 1.05 multiplicatively, we expect that our rounding does not
 252 increase ε by more than 1.05 multiplicatively.

253 In Figure 5, we plot the ratio of ε at $\delta = 10^{-6}$ for distinguishing between $\mathcal{N}(0, \hat{\sigma}^2)$ and $\mathcal{N}(X, \hat{\sigma}^2)$
 254 divided by the maximum over $i \in [T]$ of ε at $\delta = 10^{-6}$ for distinguishing between $\mathcal{N}(0, i\sigma^2)$
 255 and $\mathcal{N}(\text{Binomial}(i, q), i\sigma^2)$. We sweep over T and q , and for each q In Figure 5a (resp. Figure
 256 5b) we set σ such that distinguishing $\mathcal{N}(0, \sigma^2)$ from $\mathcal{N}(\text{Bernoulli}(q), \sigma^2)$ satisfies $(1, 10^{-6})$ -DP



(a) One iteration of DP-SGD satisfies $(1, 10^{-6})$ -DP.

(b) One iteration of DP-SGD satisfies $(2, 10^{-6})$ -DP.

Figure 5: Ratio of upper bound on ε for quadratic loss with $\alpha = 0.5$ divided by maximum ε of i iterations on a linear loss. In Figure 5a (resp. Figure 5b), for each choice of q , σ is set so 1 iteration of DP-SGD satisfies $(1, 10^{-6})$ -DP (resp $(2, 10^{-6})$ -DP).

257 (resp. $(2, 10^{-6})$ -DP). In the majority of settings, the linear loss heuristic provides a larger ε than
 258 the quadratic with $\alpha = 1/2$, and even when the quadratic provides a larger ε , the violation is small
 259 ($\leq 3\%$). This is evidence that our heuristic is still a good approximation for many convex losses.

260 4.3 Pathological Example

261 If we allow the regularizer r to be arbitrary – in particular, not even requiring continuity – then the
 262 gradient can also be arbitrary. This flexibility allows us to construct a counterexample such that the
 263 standard composition-based analysis of DP-SGD with `intermediate_iterates` is close to tight.

264 Specifically, choose the regularizer so that the update $\mathbf{m}' = \mathbf{m} - \eta \nabla_{\mathbf{m}} r(\mathbf{m})$ does the following:
 265 $\mathbf{m}'_1 = 0$ and, for $i \in [d-1]$, $\mathbf{m}'_{i+1} = v \cdot \mathbf{m}_i$. Here $v > 1$ is a large constant. We chose the loss
 266 so that, for our canary x_1 , we have $\nabla_{\mathbf{m}} \ell(\mathbf{m}, x_1) = (1, 0, 0, \dots, 0)$ and, for all other examples x_i
 267 ($i \in \{2, 3, \dots, n\}$), we have $\nabla_{\mathbf{m}} \ell(\mathbf{m}, x_i) = \mathbf{0}$. Then the last iterate is

$$\mathbf{m}_T = (A_T + \xi_{T,1}, vA_{T-1} + v\xi_{T-1,1} + \xi_{T,2}, v^2A_{T-2} + v^2\xi_{T-2,1} + v\xi_{T-1,2} + \xi_{T,3}, \dots), \quad (7)$$

268 where $A_t \leftarrow \text{Bernoulli}(p)$ indicates whether or not the canary was sampled in the t -th iteration and
 269 $\xi_{t,i}$ denotes the i -th coordinate of the noise ξ_t added in the t -th step. Essentially, the last iterate \mathbf{m}_T
 270 contains the history of all the iterates in its coordinates. Namely, the i -th coordinate of \mathbf{m}_T gives a
 271 scaled noisy approximation to A_{T-i} :

$$v^{1-i} \mathbf{m}_{T,i} = A_{T-i} + \sum_{j=0}^{i-1} v^{j+1-i} \xi_{T-j,i-j} \sim \mathcal{N}\left(A_{T-i}, \sigma^2 \frac{1-v^{-2i}}{1-v^{-2}}\right). \quad (8)$$

272 As $v \rightarrow \infty$, the variance converges to σ^2 . In other words, if v is large, from the final iterate, we can
 273 obtain $\mathcal{N}(A_i, \sigma^2)$ for all i . This makes the standard composition-based analysis of DP-SGD tight.

274 4.4 Malicious Dataset Attack

275 The examples above rely on the regularizer having large unclipped gradients. We now construct a
 276 counterexample without a regularizer, instead using other examples to amplify the canary signal.

277 Our heuristic assumes the adversary does not have access to the intermediate iterations and that the
 278 model is linear. However, we can design a nonlinear model and specific training data to directly
 279 challenge this assumption. The attack strategy is to use the model's parameters as a sort of noisy
 280 storage, saving all iterations within them. Then with access only to the final model, an adversary

Table 1: Previous works showed that large batch sizes achieve high performing models [DBHSB22]. Using our heuristic analysis it is possible to achieve similar performance for smaller batch sizes.

Batch size	Heuristic ε	Standard ε	Accuracy	Empirical ε
4096	6.34	8	79.5%	1.7
512	7.0	12	79.1%	1.8
256	6.7	14	79.4%	1.6

281 can still examine the parameters, extract the intermediate steps, and break the assumption. Our
 282 construction introduces a data point that changes its gradient based on the number of past iterations,
 283 making it easy to identify if the point was present a given iteration of training. The rest of the
 284 data points are maliciously selected to ensure the noise added during training doesn't impact the
 285 information stored in the model's parameters. We defer the full details of the attack to Appendix C.

286 Figure 6 summarizes the results. As illustrated in the figure, this attack achieves a auditing lower
 287 bound matching the standard DP-SGD analysis even in the `last_iterate_only` setting. As a result,
 288 the attack exceeds our heuristic. However, this is a highly artificial example and it is unlikely to
 289 reflect real-world scenarios.

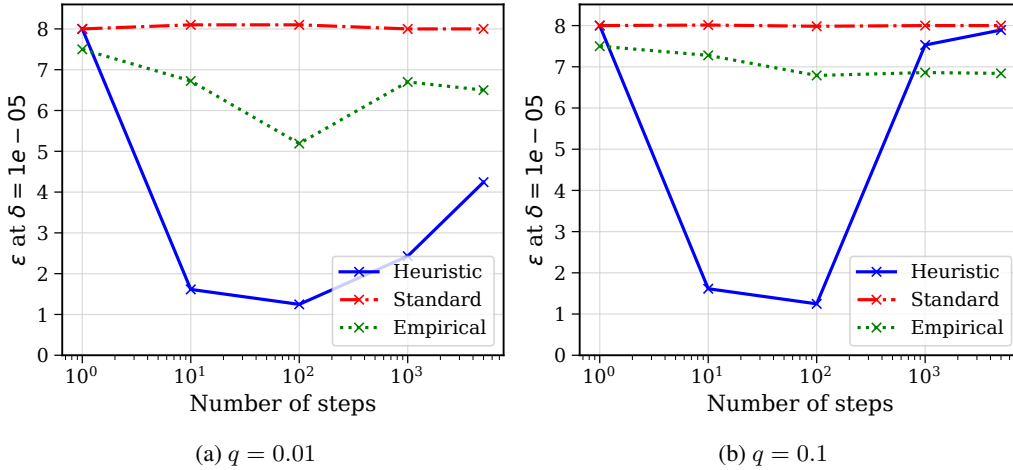


Figure 6: In this adversarial example, the attack encodes all training steps within the final model parameters, thereby violating the specific assumptions used to justify our heuristic analysis.

290 5 Discussion & Conclusion

291 Both theoretical analysis and privacy auditing are valuable for understanding privacy leakage in
 292 machine learning, but each has limitations. Theoretical analysis is inherently conservative, while
 293 auditing procedures evaluate only specific attacks, and may thus underrepresent the privacy leakage.

294 Our work introduces a novel heuristic analysis for DP-SGD that focuses on the privacy implications
 295 of releasing only the final model iterate. This approach is based in the empirical observation that
 296 linear loss functions accurately model the effectiveness of state of the art membership inference
 297 attacks. Our heuristic offers a practical and computationally efficient way to estimate privacy leakage
 298 to complement privacy auditing and the standard composition-based analysis of DP-SGD. As shown
 299 in Table 1, we trained a series of CIFAR10 models with varying batch sizes that all achieved the
 300 similar level of heuristic epsilon, albeit with different standard epsilon values. Remarkably, these
 301 models exhibited similar performance and similar empirical epsilon values.

302 We also acknowledge the limitations of our heuristic by identifying specific counterexamples where
 303 the heuristic underestimates the true privacy leakage.

304 References

- 305 [ACGMMTZ16] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar,
306 and L. Zhang. “Deep learning with differential privacy”. In: *Proceedings of*
307 *the 2016 ACM SIGSAC conference on computer and communications security*.
308 2016, pp. 308–318. URL: <https://arxiv.org/abs/1607.00133> (cit. on
309 pp. 1, 3).
- 310 [AKOOMS23] G. Andrew, P. Kairouz, S. Oh, A. Oprea, H. B. McMahan, and V. Suriyakumar.
311 “One-shot Empirical Privacy Estimation for Federated Learning”. In: *arXiv*
312 *preprint arXiv:2302.03098* (2023). URL: [https://arxiv.org/abs/2302.](https://arxiv.org/abs/2302.03098)
313 [03098](https://arxiv.org/abs/2302.03098) (cit. on p. 2).
- 314 [Ale09] K. Alex. “Learning multiple layers of features from tiny images”. In: <https://www.cs.toronto.edu/kriz/learning-features-2009-TR.pdf>
315 (2009) (cit. on pp. 5, 20).
- 317 [AT22] J. Altschuler and K. Talwar. “Privacy of noisy stochastic gradient descent:
318 More iterations without more privacy loss”. In: *Advances in Neural Information*
319 *Processing Systems* 35 (2022), pp. 3788–3800. URL: [https://arxiv.](https://arxiv.org/abs/2205.13710)
320 [org/abs/2205.13710](https://arxiv.org/abs/2205.13710) (cit. on p. 2).
- 321 [AZT24] M. Aerni, J. Zhang, and F. Tramèr. “Evaluations of Machine Learning Privacy
322 Defenses are Misleading”. In: *arXiv preprint arXiv:2404.17399* (2024). URL:
323 <https://arxiv.org/abs/2404.17399> (cit. on p. 2).
- 324 [BGDCTV18] B. Bichsel, T. Gehr, D. Drachler-Cohen, P. Tsankov, and M. Vechev. “Dp-
325 finder: Finding differential privacy violations by sampling and optimization”.
326 In: *Proceedings of the 2018 ACM SIGSAC Conference on Computer and*
327 *Communications Security*. 2018, pp. 508–524 (cit. on p. 2).
- 328 [BSA24] J. Bok, W. Su, and J. M. Altschuler. “Shifted Interpolation for Differential
329 Privacy”. In: *arXiv preprint arXiv:2403.00278* (2024). URL: [https://arxiv.](https://arxiv.org/abs/2403.00278)
330 [org/abs/2403.00278](https://arxiv.org/abs/2403.00278) (cit. on p. 2).
- 331 [BST14] R. Bassily, A. Smith, and A. Thakurta. “Private empirical risk minimization:
332 Efficient algorithms and tight error bounds”. In: *2014 IEEE 55th annual*
333 *symposium on foundations of computer science*. IEEE. 2014, pp. 464–473.
334 URL: <https://arxiv.org/abs/1405.7085> (cit. on pp. 1, 3).
- 335 [BTRKMW24] M. Bertran, S. Tang, A. Roth, M. Kearns, J. H. Morgenstern, and S. Z. Wu.
336 “Scalable membership inference attacks via quantile regression”. In: *Advances*
337 *in Neural Information Processing Systems* 36 (2024). URL: [https://arxiv.](https://arxiv.org/abs/2307.03694)
338 [org/abs/2307.03694](https://arxiv.org/abs/2307.03694) (cit. on p. 2).
- 339 [CCGST24] C. A. Choquette-Choo, A. Ganesh, T. Steinke, and A. Thakurta. *Privacy*
340 *Amplification for Matrix Mechanisms*. 2024. arXiv: [2310.15526](https://arxiv.org/abs/2310.15526) [cs.LG]
341 (cit. on p. 7).
- 342 [CCNSTT22] N. Carlini, S. Chien, M. Nasr, S. Song, A. Terzis, and F. Tramèr. “Membership
343 inference attacks from first principles”. In: *2022 IEEE Symposium on Security*
344 *and Privacy (SP)*. IEEE. 2022, pp. 1897–1914. URL: [https://arxiv.org/](https://arxiv.org/abs/2112.03570)
345 [abs/2112.03570](https://arxiv.org/abs/2112.03570) (cit. on p. 2).
- 346 [CYS21] R. Chourasia, J. Ye, and R. Shokri. “Differential privacy dynamics of langevin
347 diffusion and noisy gradient descent”. In: *Advances in Neural Information*
348 *Processing Systems* 34 (2021), pp. 14771–14781. URL: [https://arxiv.](https://arxiv.org/abs/2102.05855)
349 [org/abs/2102.05855](https://arxiv.org/abs/2102.05855) (cit. on p. 2).
- 350 [DBHSB22] S. De, L. Berrada, J. Hayes, S. L. Smith, and B. Balle. “Unlocking high-
351 accuracy differentially private image classification through scale”. In: *arXiv*
352 *preprint arXiv:2204.13650* (2022) (cit. on pp. 5, 9).
- 353 [DMNS06] C. Dwork, F. McSherry, K. Nissim, and A. Smith. “Calibrating noise to
354 sensitivity in private data analysis”. In: *Theory of Cryptography: Third Theory*
355 *of Cryptography Conference, TCC 2006, New York, NY, USA, March 4-7, 2006.*
356 *Proceedings* 3. Springer. 2006, pp. 265–284. URL: [https://www.iacr.](https://www.iacr.org/archive/tcc2006/38760266/38760266.pdf)
357 [org/archive/tcc2006/38760266/38760266.pdf](https://www.iacr.org/archive/tcc2006/38760266/38760266.pdf) (cit. on p. 1).
- 358 [DRS19] J. Dong, A. Roth, and W. J. Su. “Gaussian differential privacy”. In: *arXiv*
359 *preprint arXiv:1905.02383* (2019) (cit. on p. 13).

- 360 [DSSUV15] C. Dwork, A. Smith, T. Steinke, J. Ullman, and S. Vadhan. “Robust traceability
361 from trace amounts”. In: *2015 IEEE 56th Annual Symposium on Foundations
362 of Computer Science*. IEEE. 2015, pp. 650–669 (cit. on p. 2).
- 363 [DWWZK18] Z. Ding, Y. Wang, G. Wang, D. Zhang, and D. Kifer. “Detecting violations of
364 differential privacy”. In: *Proceedings of the 2018 ACM SIGSAC Conference
365 on Computer and Communications Security*. 2018, pp. 475–489 (cit. on p. 2).
- 366 [FMTT18] V. Feldman, I. Mironov, K. Talwar, and A. Thakurta. “Privacy amplification
367 by iteration”. In: *2018 IEEE 59th Annual Symposium on Foundations of
368 Computer Science (FOCS)*. IEEE. 2018, pp. 521–532. URL: [https://arxiv.
369 org/abs/1808.06651](https://arxiv.org/abs/1808.06651) (cit. on p. 2).
- 370 [Goo20] Google. *Differential Privacy Accounting*. [https://github.com/google/
371 differential-privacy/tree/main/python/dp_accounting](https://github.com/google/differential-privacy/tree/main/python/dp_accounting). 2020
372 (cit. on pp. 3, 4).
- 373 [HSRDTMPSNC08] N. Homer, S. Szelingner, M. Redman, D. Duggan, W. Tembe, J. Muehling,
374 J. V. Pearson, D. A. Stephan, S. F. Nelson, and D. W. Craig. “Resolving
375 individuals contributing trace amounts of DNA to highly complex mixtures
376 using high-density SNP genotyping microarrays”. In: *PLoS genetics* 4.8
377 (2008), e1000167 (cit. on p. 2).
- 378 [JUO20] M. Jagielski, J. Ullman, and A. Oprea. “Auditing differentially private ma-
379 chine learning: How private is private sgd?” In: *Advances in Neural Informa-
380 tion Processing Systems* 33 (2020), pp. 22205–22216 (cit. on p. 2).
- 381 [KJH20] A. Koskela, J. Jälkö, and A. Honkela. “Computing tight differential privacy
382 guarantees using fft”. In: *International Conference on Artificial Intelligence
383 and Statistics*. PMLR. 2020, pp. 2560–2569. URL: [https://arxiv.org/
384 abs/1906.03049](https://arxiv.org/abs/1906.03049) (cit. on p. 1).
- 385 [LF20] K. Leino and M. Fredrikson. “Stolen memories: Leveraging model memoriza-
386 tion for calibrated {White-Box} membership inference”. In: *29th USENIX se-
387 curity symposium (USENIX Security 20)*. 2020, pp. 1605–1622. URL: [https:
388 //arxiv.org/abs/1906.11798](https://arxiv.org/abs/1906.11798) (cit. on p. 2).
- 389 [Mir17] I. Mironov. “Rényi differential privacy”. In: *2017 IEEE 30th computer secu-
390 rity foundations symposium (CSF)*. IEEE. 2017, pp. 263–275. URL: [https:
391 //arxiv.org/abs/1702.07476](https://arxiv.org/abs/1702.07476) (cit. on p. 1).
- 392 [NHSBTJCT23] M. Nasr, J. Hayes, T. Steinke, B. Balle, F. Tramèr, M. Jagielski, N. Carlini,
393 and A. Terzis. “Tight Auditing of Differentially Private Machine Learning”.
394 In: *arXiv preprint arXiv:2302.07956* (2023). URL: [https://arxiv.org/
395 abs/2302.07956](https://arxiv.org/abs/2302.07956) (cit. on pp. 1, 2, 5, 6).
- 396 [NSTPC21] M. Nasr, S. Song, A. Thakurta, N. Papernot, and N. Carlini. “Adversary
397 instantiation: Lower bounds for differentially private machine learning”. In:
398 *2021 IEEE Symposium on security and privacy (SP)*. IEEE. 2021, pp. 866–
399 882. URL: <https://arxiv.org/abs/2101.04535> (cit. on pp. 1, 2, 5, 6).
- 400 [SDSOJ19] A. Sablayrolles, M. Douze, C. Schmid, Y. Ollivier, and H. Jégou. “White-
401 box vs black-box: Bayes optimal strategies for membership inference”. In:
402 *International Conference on Machine Learning*. PMLR. 2019, pp. 5558–5567.
403 URL: <https://arxiv.org/abs/1908.11229> (cit. on p. 2).
- 404 [SNJ23] T. Steinke, M. Nasr, and M. Jagielski. “Privacy auditing with one (1) training
405 run”. In: *Advances in Neural Information Processing Systems* 36 (2023). URL:
406 <https://arxiv.org/abs/2305.08846> (cit. on pp. 2, 6).
- 407 [SOJH09] S. Sankararaman, G. Obozinski, M. I. Jordan, and E. Halperin. “Genomic
408 privacy and limits of individual detection in a pool”. In: *Nature genetics* 41.9
409 (2009), pp. 965–967 (cit. on p. 2).
- 410 [SSSS17] R. Shokri, M. Stronati, C. Song, and V. Shmatikov. “Membership inference
411 attacks against machine learning models”. In: *2017 IEEE symposium on
412 security and privacy (SP)*. IEEE. 2017, pp. 3–18 (cit. on p. 2).
- 413 [Ste22] T. Steinke. “Composition of Differential Privacy & Privacy Amplification
414 by Subsampling”. In: *arXiv preprint arXiv:2210.00597* (2022). URL: [https:
415 //arxiv.org/abs/2210.00597](https://arxiv.org/abs/2210.00597) (cit. on p. 1).

416 [TTSSJC22] F. Tramer, A. Terzis, T. Steinke, S. Song, M. Jagielski, and N. Carlini. “Debug-
417 ging differential privacy: A case study for privacy auditing”. In: *arXiv preprint*
418 *arXiv:2202.12219* (2022). URL: <https://arxiv.org/abs/2202.12219>
419 (cit. on p. 2).

420 [WBKBGGG23] Y. Wen, A. Bansal, H. Kazemi, E. Borgnia, M. Goldblum, J. Geiping, and
421 T. Goldstein. “Canary in a coalmine: Better membership inference with en-
422 sembled adversarial queries”. In: *ICLR*. 2023. URL: [https://arxiv.org/](https://arxiv.org/abs/2210.10750)
423 [abs/2210.10750](https://arxiv.org/abs/2210.10750) (cit. on p. 2).

424 [YS22] J. Ye and R. Shokri. “Differentially private learning needs hidden state (or
425 much faster convergence)”. In: *Advances in Neural Information Processing*
426 *Systems* 35 (2022), pp. 703–715. URL: [https://arxiv.org/abs/2203.](https://arxiv.org/abs/2203.05363)
427 [05363](https://arxiv.org/abs/2203.05363) (cit. on p. 2).

428 [ZBWTSRPNK22] S. Zanella-Béguelin, L. Wutschitz, S. Tople, A. Salem, V. Rühle, A. Paverd,
429 M. Naseri, and B. Köpf. “Bayesian estimation of differential privacy”. In:
430 *arXiv preprint arXiv:2206.05199* (2022) (cit. on p. 2).

431 [ZK16] S. Zagoruyko and N. Komodakis. “Wide residual networks”. In: *arXiv preprint*
432 *arXiv:1605.07146* (2016) (cit. on pp. 5, 20).

433 [ZLS23] S. Zarifzadeh, P. C.-J. M. Liu, and R. Shokri. “Low-Cost High-Power Mem-
434 bership Inference by Boosting Relativity”. In: (2023). URL: [https://arxiv.](https://arxiv.org/abs/2312.03262)
435 [org/abs/2312.03262](https://arxiv.org/abs/2312.03262) (cit. on p. 2).

436 A Proof of Theorem 1

437 *Proof.* Let x_{i^*} be the canary, let D be the dataset with the canary and D' be the dataset without the
438 canary. Since ℓ and r are linear, wlog we can assume $r = 0$ and $\nabla_{\mathbf{m}_{t-1}} \ell(\mathbf{m}_{t-1}, x_i) = \mathbf{v}_i$ for some
439 set of vectors $\{\mathbf{v}_i\}$, such that $\|\mathbf{v}_i\|_2 \leq 1$. We can also assume wlog $\|\mathbf{v}_{i^*}\| = 1$ since, if $\|\mathbf{v}_{i^*}\| < 1$,
440 the final privacy guarantee we show only improves.

441 We have the following recursion for \mathbf{m}_t :

$$\mathbf{m}_t = \mathbf{m}_{t-1} - \eta \left(\sum_{i \in B_t} \mathbf{v}_i + \xi_t \right), \quad \xi_t \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2 I_d).$$

442 Unrolling the recursion:

$$\mathbf{m}_t = \mathbf{m}_0 - \eta \left[\sum_{t \in [T]} \sum_{i \in B_t} \mathbf{v}_i + \xi \right], \quad \xi \sim N(0, T\sigma^2 I_d).$$

443 By the post-processing property of DP, we can assume that in addition to the final model \mathbf{m}_T , we
444 release \mathbf{m}_0 and $\{B_t \setminus \{x_{i^*}\}\}_{t \in [T]}$, that is we release all examples that were sampled in each batch
445 except for the canary. The following f is a bijection, computable by an adversary using the released
446 information:

$$f(\mathbf{m}_T) := - \left[\frac{\mathbf{m}_T - \mathbf{m}_0}{\eta} - \sum_{t \in [T]} \sum_{i \in B_t \setminus \{x_{i^*}\}} \mathbf{v}_i \right]$$

447 Since f is a bijection, distinguishing \mathbf{m}_T sampled using D and D' is equivalent to distinguishing
448 $f(\mathbf{m}_T)$ instead. Now we have $f(\mathbf{m}_T) = \mathcal{N}(0, T\sigma^2 I_d)$ for D' , and $f(\mathbf{m}_T) = \mathcal{N}(0, T\sigma^2 I_d) + k\mathbf{v}_{i^*}$,
449 $k \sim \text{Binomial}(T, q)$. For any vector \mathbf{u} orthogonal to \mathbf{v}_{i^*} , by isotropy of the Gaussian distribution the
450 distribution of $\langle f(\mathbf{m}_T), \mathbf{u} \rangle$ is the same for both D and D' and independent of $\langle f(\mathbf{m}_T), \mathbf{v}_{i^*} \rangle$, hence
451 distinguishing $f(\mathbf{m}_T)$ given D and D' is the same as distinguishing $\langle f(\mathbf{m}_T), \mathbf{v}_{i^*} \rangle$ given D and D' .
452 Finally, the distribution of $\langle f(\mathbf{m}_T), \mathbf{v}_{i^*} \rangle$ is exactly P for D and exactly Q for D' . By post-processing,
453 this gives the theorem.

454 We can also see that the function $\delta_{T,q,\sigma}$ is tight (i.e., even if we do not release $B_t \setminus \{x_{i^*}\}$), by
 455 considering the 1-dimensional setting, where $\mathbf{v}_i = 0$ for $i \neq i^*$ and $\mathbf{v}_{i^*} = -1, \eta = 1, \mathbf{m}_0 = 0$. Then,
 456 the distribution of \mathbf{m}_T given D is exactly P , and given D' is exactly Q . \square

457 A.1 Computing δ from ε

458 Here, we give an efficiently computable expression for the function $\delta_{T,q,\sigma}(\varepsilon)$. Using P, Q as in
 459 Theorem 1, let $f(y)$ be the privacy loss for the output y :

$$\begin{aligned} f(y) &= \log \left(\frac{P(y)}{Q(y)} \right) = \log \left(\sum_{k=0}^T \binom{T}{k} q^k (1-q)^{n-k} \frac{\exp(-(y-k)^2/2T\sigma^2)}{\exp(-y^2/2T\sigma^2)} \right) \\ &= \log \left(\sum_{k=0}^T \binom{T}{k} q^k (1-q)^k \exp \left(\frac{2ky - k^2}{2T\sigma^2} \right) \right). \end{aligned}$$

460 Then for any ε , using the fact that $S = \{y : f(y) \geq \varepsilon\}$ maximizes $P(S) - e^\varepsilon Q(S)$, we have:

$$\begin{aligned} H_{e^\varepsilon}(P, Q) &= P(\{y : f(y) \geq \varepsilon\}) - e^\varepsilon Q(\{y : f(y) \geq \varepsilon\}) \\ &= P(\{y : y \geq f^{-1}(\varepsilon)\}) - e^\varepsilon Q(\{y : y \geq f^{-1}(\varepsilon)\}) \\ &= \sum_{k=0}^T \binom{T}{k} q^k (1-q)^k \Pr[\mathcal{N}(k, T\sigma^2) \geq f^{-1}(\varepsilon)] - e^\varepsilon \Pr[\mathcal{N}(0, T\sigma^2) \geq f^{-1}(\varepsilon)]. \end{aligned}$$

461 Similarly, $S = \{y : f(y) \leq -\varepsilon\}$ maximizes $Q(S) - e^\varepsilon P(S)$ so we have:

$$\begin{aligned} H_{e^\varepsilon}(Q, P) &= Q(\{y : f(y) \leq -\varepsilon\}) - e^\varepsilon P(\{y : f(y) \leq -\varepsilon\}) \\ &= Q(\{y : y \leq f^{-1}(-\varepsilon)\}) - e^\varepsilon P(\{y : y \leq f^{-1}(-\varepsilon)\}) \\ &= \Pr[\mathcal{N}(0, T\sigma^2) \leq f^{-1}(-\varepsilon)] - e^\varepsilon \sum_{k=0}^T \binom{T}{k} q^k (1-q)^k \Pr[\mathcal{N}(k, T\sigma^2) \leq f^{-1}(-\varepsilon)]. \end{aligned}$$

462 These expressions can be evaluated efficiently. Since f is monotone, it can be inverted via binary
 463 search. We can also use binary search to evaluate ε as a function of δ .

464 B Linear Worst Case for Full Batch Setting

465 It turns out that in the full-batch setting, the worst-case analyses of DP-GD with
 466 `intermediate_iterates` and with `last_iterate_only` are the same. This phenomenon arises
 467 because there is no subsampling (because $q = 1$ in Algorithm 1) and thus the algorithm is “just”
 468 the Gaussian mechanism. Intuitively, DP-GD with `intermediate_iterates` corresponds to T
 469 calls to the Gaussian mechanism with noise multiplier σ , while DP-GD with `last_iterate_only`
 470 corresponds to one call to the Gaussian mechanism with noise multiplier σ/\sqrt{T} ; these are equivalent
 471 by the properties of the Gaussian distribution.

472 We can formalize this using the language of Gaussian DP [DRS19]: DP-GD (Algorithm 1 with
 473 $q = 1$) satisfies \sqrt{T}/σ -GDP. (Each iteration satisfies $1/\sigma$ -GDP and adaptive composition implies
 474 the overall guarantee.) This means that the privacy loss is exactly dominated by that of the Gaussian
 475 mechanism with noise multiplier σ/\sqrt{T} . Linear losses give an example such that DP-GD with
 476 `last_iterate_only` has exactly this privacy loss, since the final iterate reveals the sum of all the
 477 noisy gradient estimates. The worst-case privacy of DP-GD with `intermediate_iterates` is no
 478 worse than that of DP-GD with `last_iterate_only`. The reverse is also true (by postprocessing).

479 In more detail: For T iterations of (full-batch) DP-GD on a linear losses, if the losses are (wlog)
 480 1-Lipschitz and we add noise $\mathcal{N}(0, \frac{\sigma^2}{n^2} \cdot I)$ to the gradient in every round, distinguishing the last

481 iterate of DP-SGD on adjacent databases is equivalent to distinguishing $\mathcal{N}(0, T\sigma^2)$ and $\mathcal{N}(T, T\sigma^2)$.
 482 This can be seen as a special case of Theorem 1 for $p = 1$, so we do not give a detailed argument
 483 here.

484 If instead we are given every iteration \mathbf{m}_t , for any 1-Lipschitz loss, distinguishing the joint distribu-
 485 tions of \mathbf{m}_t given \mathbf{m}_{t-1} on adjacent databases is equivalent to distinguishing $\mathcal{N}(0, \sigma^2)$ and $\mathcal{N}(1, \sigma^2)$.
 486 In turn, distinguishing the distribution of all iterates on adjacent databases is equivalent to distin-
 487 guishing $\mathcal{N}(\mathbf{0}^T, \sigma^2 I_T)$ and $\mathcal{N}(\mathbf{1}^T, \sigma^2 I_T)$, where $\mathbf{0}^T$ and $\mathbf{1}^T$ are the all-zeros and all-ones vectors in
 488 \mathbb{R}^T . Because the Gaussian distribution is isotropic, distinguishing $\mathcal{N}(\mathbf{0}^T, \sigma^2 I_T)$ and $\mathcal{N}(\mathbf{1}^T, \sigma^2 I_T)$ is
 489 equivalent to distinguishing $\langle \mathbf{x}, \mathbf{1}^T \rangle$ where $\mathbf{x} \sim \mathcal{N}(\mathbf{0}^T, \sigma^2 I_T)$ and $\langle \mathbf{x}, \mathbf{1}^T \rangle$ where $\mathbf{x} \sim \mathcal{N}(\mathbf{1}^T, \sigma^2 I_T)$.
 490 These distributions are $\mathcal{N}(0, T\sigma^2)$ and $\mathcal{N}(T, T\sigma^2)$, the exact pair of distributions we reduced to for
 491 last-iterate analysis of linear losses.

492 C Malicious Dataset Attack Details

493 Algorithms 2, 3, and 4 summarizes the construction for the attack. The attack assume the model
 494 parameters have dimension equal to the number of iterations. It also assumes each data point can
 495 reference which iteration of training is currently happening (this can be implemented by having
 496 a single model parameter which increments in each step, independently of the training examples,
 497 without impacting the privacy of the training process). Then we build our two datasets D and
 498 $D' = D \cup \{x\}$ so that all points in dataset D (“repeaters”) run Algorithm 3 to compute gradients
 499 and the canary point in D' runs Algorithm 2 to compute its gradient. Our attack relies heavily on
 500 DP-SGD’s lack of assumptions on the data distribution and any specific properties of the model or
 501 gradients. Algorithm 2, which generates the canary data point, is straightforward. Its goal is to store
 502 in the model parameters whether it was present in iteration i by outputting a gradient that changes
 503 only the i -th index of the model parameters by 1 (assuming a clipping threshold of 1).

504 All other data points, the “repeaters”, are present in both datasets (D and D'), and have three tasks:

- 505 • Cancel out any noise added to the model parameters at an index larger than the current
 506 iteration. At iteration i , their gradients for parameters from index i onward will be the same
 507 as the current value of the model parameter, scaled by the batch size and the learning rate to
 508 ensure this parameter value will be 0 after the update.
- 509 • Evaluate whether the canary point was present in the previous iteration by comparing
 510 the model parameter at index $i - 1$ with a threshold, and rewrite the value of that model
 511 parameter to a large value if the canary was present.
- 512 • Ensure that all previous decisions are not overwritten by noise by continuing to rewrite them
 513 with a large value based on their previous value.

514 To achieve all of these goals simultaneously, we require that the batch size is large enough that the
 515 repeaters’ updates are not clipped.

516 Finally Algorithm 4 runs DP-SGD, with repeater points computing gradients with Algorithm 3 and
 517 the canary point, sampled with probability p , computing its gradient using Algorithm 2. In our
 518 experiments we run Algorithm 4 100,000 times. And to evaluate if the model parameters was from
 519 dataset D or D' we run a hypothesis test on the values of the model parameters. All constants are
 520 chosen to ensure all objectives of the repeaters are satisfied.

Algorithm 2 Canary data point

```

1: function ADV( $\mathbf{x}, i$ )
2:   Initialize  $\mathbf{a}$  as a zero vector of the same dimension as  $\mathbf{x}$ 
3:   Set  $a_i \leftarrow 1$                                      ▷ Set the  $i$ -th component to 1
4:   return  $-\mathbf{a}$ 
5: end function

```

Algorithm 3 Additional data points

Require: model parameters \mathbf{x} , iteration number i , batch size N , learning rate η , previous history threshold t_{past} , last iteration threshold t_{last} , history amplification value BIG_VAL

- 1: **function** REPEATERS(\mathbf{x} , i , N , η , t_{past} , t_{last} , BIG_VAL)
- 2: $\mathbf{h} \leftarrow \mathbf{x}_{0:i}$ ▷ Parameter “history” up to iteration i , not inclusive
- 3: $\mathbf{f} \leftarrow \mathbf{x}_{i:\text{end}}$ ▷ Future and current parameters, starting from iteration i
- 4: $\mathbf{f} \leftarrow -\mathbf{f}/(\eta \cdot N)$ ▷ Remove noise from last iteration
- 5: base_history $\leftarrow -\mathbf{x}_{0:i}/(\eta \cdot N)$ ▷ By default, zero out entire history
- 6: **if** length(\mathbf{h}) > 1 **then**
- 7: $\mathbf{h}_{0:i-1} \leftarrow \text{BIG_VAL}/(\eta \cdot N) \cdot (2\mathbb{1}[\mathbf{h}_{0:i-1} \geq t_{\text{past}}] - 1)$ ▷ If an old iteration is large enough, it was a canary iteration, so amplify it
- 8: **end if**
- 9: **if** length(\mathbf{h}) > 0 **then**
- 10: $\mathbf{h}_{i-1} \leftarrow \text{BIG_VAL}/(\eta \cdot N) \cdot (2\mathbb{1}[\mathbf{h}_i \geq t_{\text{last}}] - 1)$ ▷ If the last iteration is large enough, it was a canary iteration, so amplify it
- 11: **end if**
- 12: $\mathbf{h} \leftarrow \mathbf{h} + \text{base_history}$ ▷ Don’t zero out canary iterations
- 13: $\mathbf{a} \leftarrow \text{concatenate}(\mathbf{h}, \mathbf{f})$
- 14: **return** $-\mathbf{a}$
- 15: **end function**

Algorithm 4 Encoding Attacking

Require: add-diff, whether to add the canary, batch size N , sampling rate p , learning rate (η), iteration count/parameter count D

- 1: **function** RUN_DPSGD(add-diff)
- 2: $C \leftarrow 1$
- 3: Initialize model $\mathbf{m} \leftarrow \mathbf{0}$ of dimension D
- 4: **for** $i = 0$ to D **do**
- 5: Generate a uniform random value $q \in [0, 1]$
- 6: $\mathbf{r} \leftarrow \text{repeaters}(\mathbf{m}, i)$
- 7: Compute norm $c \leftarrow \|\mathbf{r}\|$
- 8: **if** $c > 0$ **then**
- 9: Normalize $\mathbf{r} \leftarrow \mathbf{r}/\max(c, C)$
- 10: **end if**
- 11: Adjusted vector $\mathbf{z} \leftarrow \mathbf{r} \times N$
- 12: Verify condition on \mathbf{m}_i
- 13: **if** $p \leq q$ and add-diff **then**
- 14: $\mathbf{r} \leftarrow \text{adv}(\mathbf{m}, i)$
- 15: Normalize and update \mathbf{z}
- 16: **end if**
- 17: Apply Gaussian noise to \mathbf{z}
- 18: Update model $\mathbf{m} \leftarrow \mathbf{m} - \mathbf{z} \times \eta$
- 19: **end for**
- 20: **return** \mathbf{m}
- 21: **end function**

521 **NeurIPS Paper Checklist**

522 **1. Claims**

523 Question: Do the main claims made in the abstract and introduction accurately reflect the
524 paper's contributions and scope?

525 Answer: [\[Yes\]](#)

526 Justification: The abstract and introduction state what we do and then the following sections
527 and the appendix provide details.

528 Guidelines:

- 529 • The answer NA means that the abstract and introduction do not include the claims
530 made in the paper.
- 531 • The abstract and/or introduction should clearly state the claims made, including the
532 contributions made in the paper and important assumptions and limitations. A No or
533 NA answer to this question will not be perceived well by the reviewers.
- 534 • The claims made should match theoretical and experimental results, and reflect how
535 much the results can be expected to generalize to other settings.
- 536 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
537 are not attained by the paper.

538 **2. Limitations**

539 Question: Does the paper discuss the limitations of the work performed by the authors?

540 Answer: [\[Yes\]](#)

541 Justification: Section 4 discusses the limitations.

542 Guidelines:

- 543 • The answer NA means that the paper has no limitation while the answer No means that
544 the paper has limitations, but those are not discussed in the paper.
- 545 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 546 • The paper should point out any strong assumptions and how robust the results are to
547 violations of these assumptions (e.g., independence assumptions, noiseless settings,
548 model well-specification, asymptotic approximations only holding locally). The authors
549 should reflect on how these assumptions might be violated in practice and what the
550 implications would be.
- 551 • The authors should reflect on the scope of the claims made, e.g., if the approach was
552 only tested on a few datasets or with a few runs. In general, empirical results often
553 depend on implicit assumptions, which should be articulated.
- 554 • The authors should reflect on the factors that influence the performance of the approach.
555 For example, a facial recognition algorithm may perform poorly when image resolution
556 is low or images are taken in low lighting. Or a speech-to-text system might not be
557 used reliably to provide closed captions for online lectures because it fails to handle
558 technical jargon.
- 559 • The authors should discuss the computational efficiency of the proposed algorithms
560 and how they scale with dataset size.
- 561 • If applicable, the authors should discuss possible limitations of their approach to
562 address problems of privacy and fairness.
- 563 • While the authors might fear that complete honesty about limitations might be used by
564 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
565 limitations that aren't acknowledged in the paper. The authors should use their best
566 judgment and recognize that individual actions in favor of transparency play an impor-
567 tant role in developing norms that preserve the integrity of the community. Reviewers
568 will be specifically instructed to not penalize honesty concerning limitations.

569 **3. Theory Assumptions and Proofs**

570 Question: For each theoretical result, does the paper provide the full set of assumptions and
571 a complete (and correct) proof?

572 Answer: [\[Yes\]](#)

573 Justification: Theorem 1 is proved in Appendix A.

574 Guidelines:

- 575 • The answer NA means that the paper does not include theoretical results.
- 576 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
- 577 referenced.
- 578 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 579 • The proofs can either appear in the main paper or the supplemental material, but if
- 580 they appear in the supplemental material, the authors are encouraged to provide a short
- 581 proof sketch to provide intuition.
- 582 • Inversely, any informal proof provided in the core of the paper should be complemented
- 583 by formal proofs provided in appendix or supplemental material.
- 584 • Theorems and Lemmas that the proof relies upon should be properly referenced.

585 4. Experimental Result Reproducibility

586 Question: Does the paper fully disclose all the information needed to reproduce the main ex-

587 perimental results of the paper to the extent that it affects the main claims and/or conclusions

588 of the paper (regardless of whether the code and data are provided or not)?

589 Answer: [Yes]

590 Justification: The setup is described for each experiment we conduct and we reference prior

591 work that these build on.

592 Guidelines:

- 593 • The answer NA means that the paper does not include experiments.
- 594 • If the paper includes experiments, a No answer to this question will not be perceived
- 595 well by the reviewers: Making the paper reproducible is important, regardless of
- 596 whether the code and data are provided or not.
- 597 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 598 to make their results reproducible or verifiable.
- 599 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 600 For example, if the contribution is a novel architecture, describing the architecture fully
- 601 might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 602 be necessary to either make it possible for others to replicate the model with the same
- 603 dataset, or provide access to the model. In general, releasing code and data is often
- 604 one good way to accomplish this, but reproducibility can also be provided via detailed
- 605 instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 606 of a large language model), releasing of a model checkpoint, or other means that are
- 607 appropriate to the research performed.
- 608 • While NeurIPS does not require releasing code, the conference does require all submis-
- 609 sions to provide some reasonable avenue for reproducibility, which may depend on the
- 610 nature of the contribution. For example
- 611 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
- 612 to reproduce that algorithm.
- 613 (b) If the contribution is primarily a new model architecture, the paper should describe
- 614 the architecture clearly and fully.
- 615 (c) If the contribution is a new model (e.g., a large language model), then there should
- 616 either be a way to access this model for reproducing the results or a way to reproduce
- 617 the model (e.g., with an open-source dataset or instructions for how to construct
- 618 the dataset).
- 619 (d) We recognize that reproducibility may be tricky in some cases, in which case
- 620 authors are welcome to describe the particular way they provide for reproducibility.
- 621 In the case of closed-source models, it may be that access to the model is limited in
- 622 some way (e.g., to registered users), but it should be possible for other researchers
- 623 to have some path to reproducing or verifying the results.

624 5. Open access to data and code

625 Question: Does the paper provide open access to the data and code, with sufficient instruc-

626 tions to faithfully reproduce the main experimental results, as described in supplemental

627 material?

628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680

Answer: [No]

Justification: We intend to release the code eventually, but we are not able to do so at the moment; we refrain from providing a detailed reason, as this could violate anonymity.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: The setup is described for each experiment we conduct and we reference prior work that these build on. We use the standard CIFAR10 dataset for deep learning experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: The auditing results present a lower bound which can be viewed as a one-sided confidence interval. For the other results the numbers are computed non-statistically (i.e. by numerically evaluating a formula); the only potential error here is due to numerical precision.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).

- 681 • The method for calculating the error bars should be explained (closed form formula,
682 call to a library function, bootstrap, etc.)
- 683 • The assumptions made should be given (e.g., Normally distributed errors).
- 684 • It should be clear whether the error bar is the standard deviation or the standard error
685 of the mean.
- 686 • It is OK to report 1-sigma error bars, but one should state it. The authors should
687 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
688 of Normality of errors is not verified.
- 689 • For asymmetric distributions, the authors should be careful not to show in tables or
690 figures symmetric error bars that would yield results that are out of range (e.g. negative
691 error rates).
- 692 • If error bars are reported in tables or plots, The authors should explain in the text how
693 they were calculated and reference the corresponding figures or tables in the text.

694 8. Experiments Compute Resources

695 Question: For each experiment, does the paper provide sufficient information on the com-
696 puter resources (type of compute workers, memory, time of execution) needed to reproduce
697 the experiments?

698 Answer: [Yes]

699 Justification: We used A2-megagpu-16g machines from Google cloud which have 16 Nvidia
700 A100 40GB GPUs to run the experiments in this paper. Overall we used around 33,000
701 hours of GPU to run all of the experiments in the paper.

702 Guidelines:

- 703 • The answer NA means that the paper does not include experiments.
- 704 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
705 or cloud provider, including relevant memory and storage.
- 706 • The paper should provide the amount of compute required for each of the individual
707 experimental runs as well as estimate the total compute.
- 708 • The paper should disclose whether the full research project required more compute
709 than the experiments reported in the paper (e.g., preliminary or failed experiments that
710 didn't make it into the paper).

711 9. Code Of Ethics

712 Question: Does the research conducted in the paper conform, in every respect, with the
713 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

714 Answer: [Yes]

715 Justification: No human subjects or sensitive data were used.

716 Guidelines:

- 717 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 718 • If the authors answer No, they should explain the special circumstances that require a
719 deviation from the Code of Ethics.
- 720 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
721 eration due to laws or regulations in their jurisdiction).

722 10. Broader Impacts

723 Question: Does the paper discuss both potential positive societal impacts and negative
724 societal impacts of the work performed?

725 Answer: [NA]

726 Justification: This work is primarily theoretical. While it is possible that downstream uses
727 of our work could be societally impactful, the precise consequences are difficult to foresee.
728 The considerations are similar to any other paper on private machine learning.

729 Guidelines:

- 730 • The answer NA means that there is no societal impact of the work performed.

- 731 • If the authors answer NA or No, they should explain why their work has no societal
732 impact or why the paper does not address societal impact.
- 733 • Examples of negative societal impacts include potential malicious or unintended uses
734 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
735 (e.g., deployment of technologies that could make decisions that unfairly impact specific
736 groups), privacy considerations, and security considerations.
- 737 • The conference expects that many papers will be foundational research and not tied
738 to particular applications, let alone deployments. However, if there is a direct path to
739 any negative applications, the authors should point it out. For example, it is legitimate
740 to point out that an improvement in the quality of generative models could be used to
741 generate deepfakes for disinformation. On the other hand, it is not needed to point out
742 that a generic algorithm for optimizing neural networks could enable people to train
743 models that generate Deepfakes faster.
- 744 • The authors should consider possible harms that could arise when the technology is
745 being used as intended and functioning correctly, harms that could arise when the
746 technology is being used as intended but gives incorrect results, and harms following
747 from (intentional or unintentional) misuse of the technology.
- 748 • If there are negative societal impacts, the authors could also discuss possible mitigation
749 strategies (e.g., gated release of models, providing defenses in addition to attacks,
750 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
751 feedback over time, improving the efficiency and accessibility of ML).

752 11. Safeguards

753 Question: Does the paper describe safeguards that have been put in place for responsible
754 release of data or models that have a high risk for misuse (e.g., pretrained language models,
755 image generators, or scraped datasets)?

756 Answer: [NA]

757 Justification: Our paper uses standard datasets (CIFAR10) and standard models (WideRes-
758 Net).

759 Guidelines:

- 760 • The answer NA means that the paper poses no such risks.
- 761 • Released models that have a high risk for misuse or dual-use should be released with
762 necessary safeguards to allow for controlled use of the model, for example by requiring
763 that users adhere to usage guidelines or restrictions to access the model or implementing
764 safety filters.
- 765 • Datasets that have been scraped from the Internet could pose safety risks. The authors
766 should describe how they avoided releasing unsafe images.
- 767 • We recognize that providing effective safeguards is challenging, and many papers do
768 not require this, but we encourage authors to take this into account and make a best
769 faith effort.

770 12. Licenses for existing assets

771 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
772 the paper, properly credited and are the license and terms of use explicitly mentioned and
773 properly respected?

774 Answer: [Yes]

775 Justification: We use CIFAR10 [Ale09] and a WideResNet [ZK16].

776 Guidelines:

- 777 • The answer NA means that the paper does not use existing assets.
- 778 • The authors should cite the original paper that produced the code package or dataset.
- 779 • The authors should state which version of the asset is used and, if possible, include a
780 URL.
- 781 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 782 • For scraped data from a particular source (e.g., website), the copyright and terms of
783 service of that source should be provided.

- 784
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- 785
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- 786
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.
- 787
- 788
- 789
- 790
- 791

792 13. **New Assets**

793 Question: Are new assets introduced in the paper well documented and is the documentation
794 provided alongside the assets?

795 Answer: [NA]

796 Justification: Our main contribution is a heuristic privacy analysis. This is fully described in
797 the paper and can be computed using existing open-source libraries.

798 Guidelines:

- The answer NA means that the paper does not release new assets.
 - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
 - The paper should discuss whether and how consent was obtained from people whose asset is used.
 - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.
- 799
- 800
- 801
- 802
- 803
- 804
- 805
- 806

807 14. **Crowdsourcing and Research with Human Subjects**

808 Question: For crowdsourcing experiments and research with human subjects, does the paper
809 include the full text of instructions given to participants and screenshots, if applicable, as
810 well as details about compensation (if any)?

811 Answer: [NA]

812 Justification: The paper does not involve crowdsourcing nor research with human subjects.

813 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
 - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
 - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.
- 814
- 815
- 816
- 817
- 818
- 819
- 820
- 821

822 15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 823 Subjects**

824 Question: Does the paper describe potential risks incurred by study participants, whether
825 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
826 approvals (or an equivalent approval/review based on the requirements of your country or
827 institution) were obtained?

828 Answer: [NA]

829 Justification: The paper does not involve crowdsourcing nor research with human subjects.

830 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
 - Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- 831
- 832
- 833
- 834
- 835

836
837
838
839
840

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.