

---

# Why Differentially-Private Local SGD – An Analysis of Biased Synchronized-Only Iterates

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We argue to use Differentially-Private Local Stochastic Gradient Descent (DP-  
2 LSGD) in both centralized and distributed setups, and explain why DP-LSGD  
3 enjoys higher clipping efficiency and produces less clipping bias compared to clas-  
4 sic Differentially-Private Stochastic Gradient Descent (DP-SGD). For both convex  
5 and non-convex optimization, we present generic analysis on noisy synchronized-  
6 only iterates in LSGD, the building block of federated learning, and study its  
7 applications to differentially-private gradient methods with clipping-based sen-  
8 sitivity control. We point out that given the current *decompose-then-compose*  
9 framework, there is no essential gap between the privacy analysis of centralized  
10 and distributed learning, and DP-SGD is a special case of DP-LSGD. We thus build  
11 a unified framework to characterize the clipping bias via the second moment of  
12 local updates, which initiates a direction to systematically instruct DP optimization  
13 by variance reduction. We show DP-LSGD with multiple local iterations can  
14 produce more concentrated local updates and then enables a more efficient exploita-  
15 tion of the clipping budget with a better utility-privacy tradeoff. In addition, we  
16 prove that DP-LSGD can converge faster to a small neighborhood of global/local  
17 optimum compared to regular DP-SGD. Thorough experiments on practical deep  
18 learning tasks are provided to support our developed theory.

## 19 1 Introduction

20 Local Stochastic Gradient Descent (LSGD) [1, 2] and (Local/Client-Level) Differential Privacy (DP)  
21 [3, 4, 5] are two popular methods to address the issues of communication efficiency and data privacy,  
22 respectively. Rooted in the *FedAvg* framework first proposed in [6], instead of communicating and  
23 synchronizing on the local updates from each user at each iteration, LSGD [1] randomly samples  
24 participants to perform gradient descent on their local data in parallel and only aggregates their local  
25 updates periodically. Though LSGD is a simple generalization of SGD to a distributed setup with a  
26 lower synchronization frequency, empirically it is known to produce promising performance, with  
27 regard to both communication efficiency and convergence rate [7]. When each user holds i.i.d. data,  
28 LSGD provably achieves a linear speedup in the number of users with also asymptotic improvements  
29 on the communication overhead over regular distributed SGD to produce equivalent accuracy [1, 2].

30 As for privacy preservation, DP [3, 8] provides a semantically precise way to quantify the data leakage  
31 from any processing. At a high level, DP is an input-independent guarantee which ensures that an ad-  
32 versary cannot infer the participation of an individual datapoint easily from the release. For example,  
33 the classic  $(\epsilon, \delta)$ -DP with small security parameters  $\epsilon$  and  $\delta$  implies a large Type I or Type II error for  
34 an adversarial hypothesis testing to guess whether an arbitrary individual is involved in the processing  
35 [9]. In DP research, one key problem is to determine the *sensitivity*, the worst-case influence/change  
36 on the output of the objective processing after arbitrarily replacing an individual in an input set. Only

37 with tractable sensitivity, one can then apply proper randomization/perturbation such as the Gaussian  
38 or Laplace mechanism [10] to produce required security parameters. Unfortunately, sensitivity is  
39 in general NP-hard to compute [11]. To this end, in practice, a commonly-applied alternative is the  
40 *decompose-then-compose* framework: a complicated processing is first (approximately) decomposed  
41 into several simpler (possibly adaptive) subroutines such as mean estimation, each of whose sen-  
42 sitivity is controllable. A *white-box* adversary is then assumed who can observe the intermediate  
43 computations, and an upper bound on the privacy loss is derived by the composition of the leakage  
44 from the virtual release in each step [12].

45 In the applications of machine learning, where the processing function returns a model trained on  
46 possibly sensitive data, arguably the most popular and generic DP privatization method is DP-SGD  
47 [13, 14]. As a representative of the above-mentioned decompose-then-compose framework, DP-SGD  
48 views the SGD as a sequence of adaptive gradient mean estimations. To ensure a bounded sensitivity  
49 guarantee, each per-sample gradient is clipped, usually, in  $l_2$ -norm [14] to some constant  $c$ , which is  
50 essentially a projection to an  $l_2$ -norm ball of radius  $c$ . Noise, which is determined by both the number  
51 of iterations  $T$  and the clipping threshold  $c$  (sensitivity bound), is then added to the clipped stochastic  
52 gradient in each iteration to produce satisfied DP parameters  $(\epsilon, \delta)$  under  $T$ -fold composition. A wider  
53 dimension and a longer convergence time  $T$  will consequently require a larger DP noise. Though the  
54 implementation of DP-SGD does not require any additional assumptions on either model or training  
55 data, it is notorious for heavy utility loss, especially for deep learning. Moreover, the understanding  
56 of the clipping bias from this artificial sensitivity control remains limited. In general, due to the bias,  
57 clipped SGD will not converge even without noise perturbation [15, 16].

58 Given the artificial assumption that DP-SGD releases the intermediate computations, there is no  
59 essential gap between the privacy analysis of the centralized and local SGD, except that in the  
60 distributed setup one may apply different DP metrics such as Local DP (LDP) [4] or client-level DP  
61 [5] to consider the privacy preservation for each user's local data. More interestingly, it is worth  
62 noting the connection among different problems in federated learning and DP-SGD that are essentially  
63 equivalent. First, it is not hard to see that DP-SGD is a special case of DP-LSGD. DP-SGD can  
64 be viewed as:  $n$  nodes, each holds a sample, and a virtual server collects the clipped stochastic  
65 gradient from a subset of sampled nodes *in every iteration*, and publishes a noisy gradient descent.  
66 DP-LSGD can be similarly defined where the only difference is that the server may not synchronize  
67 on each iteration, but clips and aggregates a linear combinations of local gradients, *periodically*.  
68 Thus, as a primary concern in federated learning, a smaller communication overhead in a lower  
69 synchronization/aggregation frequency would also imply less leakage and a smaller composition  
70 bound of privacy loss. On the other hand, the study on the utility loss by perturbation and artificial  
71 sensitivity control (clipping) could also be used to analyze federated learning with compressed  
72 communication [17] where there exists quantification error in broadcasted local updates. Therefore,  
73 in this paper, we aim to provide a unified analysis for both noisy LSGD and DP-LSGD/SGD to get  
74 new insights. Before we can build useful theory to capture these concerns from different perspectives,  
75 several technical challenges need to be addressed.

76 **Utility of "Synchronized/Published" Iterate Only:** Many existing convergence results [2, 18, 19,  
77 20, 21] on non-private LSGD are developed on the (weighted) average of all iterates. These include  
78 the intermediate iterates produced during the local updates from each user/node, which will not be  
79 exposed or shared. To properly characterize the effect of perturbation, a more appropriate and realistic  
80 convergence guarantee is to measure the performance of synchronized (shared) iterates only. This is  
81 also important to help understand the practical performance of LSGD as neither the server nor users  
82 have access to all intermediate computations. Such measurement is especially necessary when we  
83 apply LSGD in a private version: the utility of concern is only with respect to the released outputs,  
84 and anything assumed to be published would incur privacy loss and increase the scale of DP noise.

85 **Clipping Bias and Data Heterogeneity:** In practice, tight sensitivity of many data processing  
86 algorithms is intractable and thus a very popular but artificial control is clipping. However, clipping  
87 could also bring non-negligible bias. In general, there is no convergence guarantee for clipped SGD  
88 if we only assume the stochastic gradient is of bounded variance [15], though under more restrictive  
89 assumptions, for example, when the stochastic gradient is in a symmetric [15] or light-tailed [22]  
90 distribution, or provided generalized smoothness [23], some (near) convergence results are known. A  
91 concise characterization of such clipping bias still largely remains open, especially for deep learning.  
92 The bias is even more complicated in the more general DP-LSGD. To provide meaningful theory  
93 to instruct systematic bias reduction, we do not want to assume Lipschitz continuity or bounded

94 gradient, which may make the analysis trivial and impractical. Thus, the desired analysis essentially  
95 captures the scenario given heavy data heterogeneity, and the results should not require a bounded  
96 difference among the local updates.

97 In this paper, through tackling the above-mentioned challenges, we aim to provide useful and intuitive  
98 theory to understand practical performance of LSGD and instruct optimization with DP guarantees.  
99 In particular, we want to explain how DP-LSGD out-performs regular DP-SGD. We summarize our  
100 contributions as follows.

- 101 1. With only a mild assumption that the stochastic gradient is of bounded variance, we present  
102 the convergence analysis on the released-only iterates of LSGD under perturbation for both  
103 convex and non-convex smooth optimization in Theorem 3.1 and 3.2. In particular, for the  
104 general convex case, we show more powerful last iterate convergence, which could be of  
105 independent interest in developing generic last-iterate analysis with unbounded gradients.
- 106 2. We then generalize our results to study the utility of DP-LSGD, where DP-SGD becomes  
107 a special case. In particular, we use the incremental norm of local update (see Definition  
108 4.1) to characterize the clipping bias and show DP-LSGD has a faster convergence rate to a  
109 small neighborhood of global/local optimum as compared to DP-SGD.
- 110 3. We further show LSGD behaves as an efficient variance reduction of local update, where  
111 multiple local GDs with a small learning rate cancel out substantial sampling noise, and  
112 enable more efficient clipping compared to DP-SGD. Thorough experiments show that  
113 DP-LSGD produces a much sharpened utility-privacy tradeoff in practical deep learning.

## 114 1.1 Related Works

115 **Convergence Analysis of LSGD:** With the increasing scale of both training data and models,  
116 federated learning has become an important paradigm in modern machine learning, where LSGD and  
117 its variants form the building block. Though the idea of LSGD can be traced back to earlier works  
118 [24, 25], the theoretical convergence analysis has only been proved recently. A common strategy to  
119 show convergence is to consider a virtual average of all the intermediate iterates produced by each  
120 user, and keep track of the divergence (dissimilarity) between the virtual average and the local iterate.  
121 In the setup where each user holds i.i.d. data, Stich in [1] studied strongly-convex optimization with  
122 LSGD and showed a linear speedup in the number of users/nodes. [26] presented non-convex analysis  
123 under the Lipschitz continuity assumption where the divergence of local update is also bounded.

124 For the more general applications with heterogeneous data, [27] studied the convex case with local  
125 GD (without sampling on either users or users' local data) but still under Lipschitz continuity. [2]  
126 presented more generic and tighter analysis for LSGD without assumptions on bounded gradient for  
127 both strongly and general convex optimization. Further generalization of LSGD to the decentralized  
128 setup under arbitrary network topology was considered in [19, 28, 29]. However, many existing  
129 works [2, 19, 28] only showed the convergence rate relying on all the intermediate averages. To our  
130 knowledge, the first generic analysis for synchronized-only iterates was shown in [30]. [30] proposed  
131 Scaffold, a generalized LSGD with careful correction on the client-drift caused by data heterogeneity.  
132 Compared to existing works, in this paper, we prove more powerful last-iterate analysis for general  
133 convex optimization with clipping and perturbation for privacy. It is also worth mentioning that with  
134 a different motivation, there is another line of works also studying noisy LSGD to capture the effect  
135 of compressed local updates to further save the communication cost. But, in most existing related  
136 works [17, 31], the compression error is assumed to be independent with zero-mean. As we need to  
137 study DP-LSGD with clipped local update, which introduces bias in the local update generation, in  
138 this paper we present more involved analysis to handle such adaptive and biased perturbation.

139 **Convergence Analysis of DP-SGD and DP-LSGD:** Asymptotically, under Lipschitz continuity, DP-  
140 SGD is known to produce a tight utility-privacy tradeoff [32, 33], where no bias is produced given a  
141 clipping threshold larger than the Lipschitz constant. However, without Lipschitz continuity, practical  
142 understanding of DP-SGD remains limited. On one hand, negative examples are shown in [15, 16]  
143 where clipped-SGD in general will not converge, and in practice clipped-SGD does produce bias  
144 and has a lower convergence rate, especially in deep learning applications compared to regular SGD  
145 [16]. On the other hand, under more restrictive assumptions on the stochastic gradient distribution,  
146 clipped-SGD can be shown to (nearly) converge [15, 22, 23]. A generic characterization on the  
147 clipping bias still largely remains open. As a consequence, there is little known meaningful theory to

148 systematically instruct optimization algorithms with DP guarantees, and most existing private deep  
 149 learning works are empirical, which aim to search for the optimal model and hyperparameters for  
 150 objective training data [34, 35, 36]. As for DP-LSGD, to our knowledge the only known theoretical  
 151 result that captures the clipping bias is [16]. However, [16] still assumes globally bounded gradient  
 152 compared to bounded second moment as assumed in our results, and its main motivation is to study  
 153 the clipping effect in client-level DP. In this paper, we show more intuitive and generic analysis of  
 154 DP-LSGD for both convex and non-convex optimization, and our motivations are also very different.  
 155 We set out to provide usable quantification on the utility loss due to clipping and *we argue to apply*  
 156 *DP-LSGD both in the centralized and distributed setup*, since DP-LSGD can significantly reduce the  
 157 clipping bias with a faster convergence rate.

## 158 2 Preliminaries

159 We focus on the classic Empirical Risk Minimization (ERM) problem. Given a dataset  $\mathcal{D} =$   
 160  $\{(x_i, y_i), i = 1, 2, \dots, n\}$ , the loss function is defined as  $F(w) = \frac{1}{n} \cdot \sum_{i=1}^n f(w, x_i, y_i) = \frac{1}{n} \cdot$   
 161  $\sum_{i=1}^n f_i(w)$ . We will consider the cases where the loss function  $f_i(w) : \mathcal{W} \rightarrow \mathbb{R}^+$  is convex or  
 162 non-convex.  $w^* = \arg \min_w F(w)$  represents the global optimum. Some formal definitions about  
 163 the properties of the objective loss function are defined as follows.

164 **Definition 2.1** (Smoothness). *A function  $f$  is  $\beta$ -smooth on  $\mathcal{W}$  if the gradient  $\nabla f(w)$  is  $\beta$ -Lipschitz*  
 165 *such that for all  $w, w' \in \mathcal{W}$ ,  $\|\nabla f(w) - \nabla f(w')\| \leq \beta \|w' - w\|$ .*

166 **Definition 2.2** (Convexity and Strong Convexity). *A function  $f(w)$  is  $\lambda$ -convex on  $\mathcal{W}$  if for all*  
 167  *$w, w' \in \mathcal{W}$ ,  $\frac{\lambda}{2} \|w - w'\|^2 \leq f(w) - f(w') - \langle \nabla f(w'), w - w' \rangle$ . We call  $f(w)$  general convex if*  
 168  *$\lambda = 0$ , and  $f(w)$  is strongly convex if  $\lambda > 0$ .*

169 **Assumption 2.1** (Bounded Variance of Stochastic Gradient). *For any  $w \in \mathcal{W}$  and an index  $i$  that is*  
 170 *randomly selected from  $\{1, 2, \dots, n\}$ , there exists  $\tau > 0$  such that  $\mathbb{E}[\|\nabla F(w) - \nabla f_i(w)\|^2] \leq \tau$ .*

171 Assumption 2.1 is the only additional assumption we need for the analysis of non-private LSGD  
 172 without clipping. We formally present the non-private LSGD algorithm in Algorithm 1 which uses  
 173 non-clipped local update (3). The whole process is formed of  $T$  phases. In each phase, by  $q$ -Poisson  
 174 sampling, in expectation ( $nq$ ) many users will be selected to perform  $K$  local gradient descents  
 175 on their local data before broadcasting the local update. To match the DP-LSGD where the local  
 176 function  $f_i(w)$  held by each user may only be determined by a single datapoint, we do not consider  
 177 an additional stochastic gradient oracle on the local function in Algorithm 1, but only assume random  
 178 sampling on the user level at each phase. However, our results can be easily generalized to the  
 179 scenario with stochastic local gradient. Moreover, we assume Poisson sampling in Algorithm 1 so as  
 180 to match the setup of DP-LSGD, since given current studies on privacy amplification by sampling,  
 181 Poisson sampling can produce the tightest results [37] (and has become the most popular option in  
 182 practice [36, 38]). In the following, we introduce the definition of DP.

183 **Definition 2.3** (Differential Privacy [38]). *Given a universe  $\mathcal{X}^*$ , we say that two datasets  $X, X' \subseteq \mathcal{X}^*$*   
 184 *are adjacent, denoted as  $X \sim X'$ , if  $X = X' \cup x$  or  $X' = X \cup x$  for some additional datapoint*  
 185  *$x \in \mathcal{X}$ . A randomized algorithm  $\mathcal{M}$  is said to be  $(\epsilon, \delta)$ -differentially-private (DP) if for any pair of*  
 186 *adjacent datasets  $X, X'$  and any event set  $O$  in the output domain of  $\mathcal{M}$ , it holds that*

$$\mathbb{P}(\mathcal{M}(X) \in O) \leq e^\epsilon \cdot \mathbb{P}(\mathcal{M}(X') \in O) + \delta.$$

187 In Definition 2.3, we apply the unbounded DP definition as adopted in most existing DP-SGD works  
 188 [16, 35, 38], where the two adjacent datasets are defined to differ in one datapoint. One may also  
 189 apply the bounded DP definition [8] by defining the adjacent datasets as arbitrarily replacing a  
 190 datapoint. However, as a stronger definition, bounded DP will also face a larger sensitivity bound.

191 We can now formally describe DP-LSGD and DP-SGD. In (2) of Algorithm 1, a clipping operation  
 192 on a vector  $v$  with threshold  $c$  is defined as  $\mathcal{CP}(v, c) = v \cdot \min\{1, c/\|v\|\}$ , which ensures a bounded  
 193 sensitivity up to  $c$ . Using the clipped local update (2), by selecting  $Q^{(t)}$  to be proper DP noise,  
 194 Algorithm 1 captures DP-SGD when  $K = 1$  and DP-LSGD for general  $K \geq 1$ . DP-LSGD (SGD) is  
 195 essentially an LSGD (SGD) with clipped local update (per-sample gradient) and additional DP noise.  
 196 Running for  $T$  iterations with a total privacy budget  $(\epsilon, \delta)$ , one may select  $Q^{(t)} \sim \mathcal{N}(0, \sigma^2 \cdot \mathbf{I}_d)$   
 197 where  $\sigma = \tilde{O}(qc\sqrt{T \log(1/\delta)})/\epsilon$  by the composition bound [38]. The privacy analysis and the noise  
 198 bound are identical for both DP-LSGD and DP-SGD given the same clipping threshold  $c$ .

---

**Algorithm 1** (Differentially Private) Local SGD with Noisy (Clipped) Periodic Averaging

---

1: **Input:** A system of  $n$  workers where each holds a local loss function  $F(w) = f_i(w)$ , sampling rate  $q$ , update step size  $\eta$ , local update length  $K$  and global synchronization number  $T$ , clipping threshold  $c$ , and initialization  $\bar{w}^{(0)}$  with synchronization noise  $Q^{(1:T)}$ .

2: **for**  $t = 1, 2, \dots, T$  **do**

3:   Implement i.i.d. sampling to select an index batch  $S^{(t)} = \{[1], \dots, [B_t]\}$  from  $\{1, 2, \dots, n\}$  of size  $B_t$ .

4:   **for**  $i = 1, 2, \dots, B_t$  in parallel **do**

5:      $w_{[i]}^{(t,0)} = \bar{w}^{(t-1)}$ .

6:     **for**  $k = 1, 2, \dots, K$  **do**

7:        $w_{[i]}^{(t,k)} = w_{[i]}^{(t,k-1)} - \eta \nabla f_{[i]}(w_{[i]}^{(t,k-1)})$ . (1)

8:     **end for**

9:     Clip the local update as  $\Delta w_{[i]}^{(t)} = \mathcal{CP}(w_{[i]}^{(t,K)} - \bar{w}^{(t-1)}, c)$

10:   **end for**

11:   **if** to ensure Differential Privacy with clipping **then**

12:     
$$\bar{w}^{(t)} = \bar{w}^{(t-1)} + \frac{1}{nq} \cdot \left( \sum_{i=1}^{B_t} \Delta w_{[i]}^{(t)} \right) + Q^{(t)}$$
 (2)

13:   **else**

14:     
$$\bar{w}^{(t)} = \frac{1}{nq} \cdot \left( \sum_{i=1}^{B_t} w_{[i]}^{(t,K)} \right) + Q^{(t)}$$
 (3)

15:   **end if**

16: **end for**

17: **Output:**  $\bar{w}^{(t)}$  for  $t = 1, 2, \dots, T$ .

---

199 We want to stress again that our motivation to study DP-LSGD is not because we only focus on the  
200 federated setup, but to provide a unified analysis of the clipping bias and argue for using DP-LSGD  
201 *even in the centralized setup*. Our results are straightforwardly applicable to distributed learning with  
202 local DP [4] or client-level DP [5], where the only difference is that we may add a larger noise  $Q^{(t)}$   
203 determined by the number of local datapoints or the users involved, respectively, for these stronger  
204 DP definitions. As for the possible communication restriction where we need to add discrete noise of  
205 finite precision, one may replace the Gaussian noise by the Binomial mechanism [39].

### 206 3 Convergence of Synchronized-Only Iterate in Noisy Non-Clipped LSGD

207 In this section, we will study the convergence analysis of LSGD in Algorithm 1 using the non-clipped  
208 local update (3) for both convex and non-convex optimization.

209 **Theorem 3.1** (Last-iterate Convergence of Noisy LSGD in General Convex Optimization). *For an*  
210 *objective function  $F(w) = \frac{1}{n} \cdot \sum_{i=1}^n f_i(w)$  where  $f_i(w)$  is convex and  $\beta$ -smooth with variance-*  
211 *bounded gradient (Assumption 2.1), when  $\eta < \min\{\frac{\beta}{\sqrt{24K}}, \frac{1}{\beta}, \frac{1}{2\beta+3K\beta/(nq)}\}$ ,  $\log(TK) \geq 2$ , and*  
212  *$Q^{(t)}$  is an independent noise such that  $\mathbb{E}[Q^{(t)}] = 0$  and  $\mathbb{E}[\|Q^{(t)}\|^2] \leq \bar{Q}$ , for some parameter  $\bar{Q}$  for*  
213  *$t = 1, 2, \dots, T$ , Algorithm 1 with (3) ensures*

$$\begin{aligned} \mathbb{E}[F(\bar{w}^{(T)})] &\leq \left( \frac{\|\bar{w}^{(0)} - w^*\|^2}{\eta(TK+1)} + \log(TK+1)(6\eta\tau/(nq) + 8K^2\beta\tau\eta^2 + \bar{Q}/\eta) \right. \\ &\quad \left. + 5\eta\beta^2(\log(TK)+1)(\|\bar{w}^{(0)} - w^*\|^2 + T(8\beta\eta^3K^3\tau + \frac{12K^3\beta^2\eta^4\tau + 3K^2\eta^2\tau}{nq} + \bar{Q})) \right) \\ &= \tilde{O}\left( \frac{\|\bar{w}^{(0)} - w^*\|^2}{\sqrt{TK}} + \frac{\tau}{\sqrt{TK}nq} + \frac{K\tau}{T} + \sqrt{TK}\bar{Q} \right), \text{ if } \eta = O(1/\sqrt{TK}). \end{aligned}$$

214

215 The proof can be found in Appendix [A](#). To prove Theorem [3.1](#) with a careful analysis on  $\|\bar{w}^{(t)} - w^*\|^2$ ,  
 216 we develop a new last-iterate analysis framework, different from existing works [\[40, 41, 42\]](#) which  
 217 must count on the assumption of bounded gradient. In Theorem [3.1](#), we need to assume the noise  
 218  $Q$  to be independent and of zero-mean. Because we do not assume Lipschitz continuity of  $F(w)$ ,  
 219 we cannot provide a meaningful upper bound of the deviation between  $F(w)$  and  $F(w + Q)$  for  
 220 arbitrary  $w$  and  $Q$  in general. However, provided the Lipschitz assumption, Theorem [3.1](#) can be  
 221 easily generalized to handle biased perturbation. In Section [4](#), with an additional assumption on the  
 222 similarity of the local functions (Assumption [4.2](#)), we will show how to handle the clipping bias as a  
 223 special biased noise. When there is no noise  $Q = 0$ , provided that  $K = O(T^{1/3}/(nq)^{2/3})$ , we show  
 224 LSGD achieves  $\tilde{O}(\frac{\|\bar{w}^{(0)} - w^*\|^2 + \tau/(nq)^{2/3}}{\sqrt{TK}})$  last-iterate convergence in general-convex optimization.

225 We now study the non-convex scenario.

226 **Theorem 3.2** (Synchronized-only Iterate Convergence of Noisy LSGD in Non-convex Optimization).  
 227 For an arbitrary objective function  $F(w) = \frac{1}{n} \cdot \sum_{i=1}^n f_i(w)$ , where  $f_i(w)$  is  $\beta$ -smooth and satisfies  
 228 Assumption [2.1](#) and for arbitrary perturbation (not necessarily independent or of zero mean) where  
 229  $\mathbb{E}[\|Q^{(t)}\|^2] \leq \bar{Q}$ , when  $\eta < \min\{\frac{\beta}{\sqrt{24K}}, \frac{1}{4\beta K}\}$ , Algorithm [1](#) with [\(3\)](#) ensures that

$$\begin{aligned} \mathbb{E}\left[\frac{\sum_{t=1}^T \|\nabla F(\bar{w}^{(t-1)})\|^2}{T}\right] &\leq \frac{4F(\bar{w}^{(0)})}{TK\eta} + \frac{16\eta^2\tau\beta^2K^2}{nq} + \frac{4(1+\beta\eta)\sum_{t=1}^T \mathbb{E}[\|Q_i^{(t)}\|^2]}{\eta^2KT} \\ &= O\left(\frac{\tau^{1/3}}{T^{2/3}(nq)^{1/3}} + \frac{T^{2/3}\tau^{2/3}K\bar{Q}}{(nq)^{2/3}}\right), \end{aligned} \quad (4)$$

when we select  $\eta = O(\frac{(nq)^{1/3}}{T^{1/3}K\tau^{1/3}})$ . In particular, when  $Q^{(t)}$  is independent and  $\mathbb{E}[Q^{(t)}] = 0$ , and  
 $\eta = \Theta(1/K)$ , then

$$\mathbb{E}\left[\frac{\sum_{t=1}^T \|\nabla F(\bar{w}^{(t-1)})\|^2}{T}\right] \leq O\left(\frac{F(\bar{w}^{(0)})}{\eta TK} + \tau + \frac{\sum_{t=1}^T \beta \mathbb{E}[\|Q^{(t)}\|^2]}{\eta TK}\right) = O\left(\frac{1}{T} + \tau + \bar{Q}\right).$$

230 The proof can be found in Appendix [B](#). In Theorem [3.2](#), we provide an analysis on the effect of generic  
 231 perturbation, which can also be used to capture the clipping bias in DP-LSGD. When there is no  
 232 perturbation, Theorem [3.2](#) has two implications. First, we show to ensure  $\min \mathbb{E}[\|\nabla F(\bar{w}^{(t)})\|^2] \leq \kappa$ ,  
 233 we need  $T = O(\frac{\sqrt{\tau/(nq)}}{\kappa^{3/2}})$ , which is tighter than the state-of-the-art results  $O(\frac{\tau/(nq)}{\kappa^2} + \frac{\sqrt{\tau}}{\kappa^{3/2}})$  in  
 234 [\[30\]](#). Second, compared to  $O(1/T^{2/3})$ , we also show that LSGD can converge faster in  $O(1/T)$   
 235 to a  $\tau$ -neighborhood of a saddle point. This is helpful to understand the practical performance of  
 236 DP-LSGD with bias, as discussed in Section [4.2](#).  
 237

238 As a final remark, we want to mention it is possible to improve the convergence rate from  $O(1/T^{2/3})$   
 239 to  $O(1/T)$  via careful variance reduction or error feedback mechanism, such as Scaffold [\[30\]](#) or  
 240 FedLin [\[43\]](#). However, the proper implementation of those advanced tricks in DP-LSGD with  
 241 additional sensitivity control is not clear. As a first step to systematically study the generic clipping  
 242 bias, in this paper we only focus on the regular LSGD. We will explain and discuss possible  
 243 generalizations in Section [6](#).

## 244 4 Utility and Clipping Bias of DP-LSGD and DP-SGD

245 In this section, we move to study DP-LSGD with clipped local update [\(2\)](#) in Algorithm [1](#). To have  
 246 a clear comparison with DP-SGD, we still consider the centralized setup and  $F(w) = 1/n \cdot f_i(w)$   
 247 where each local function  $f_i(w)$  is determined by a single sample. To capture the clipping bias, we  
 248 need to introduce a new term, termed *incremental norm*.

249 **Definition 4.1** (Incremental Norm). Consider applying the private and clipping version of Algorithm [1](#)  
 250 with [\(2\)](#) on  $F(w) = \sum_{i=1}^n f_i(w)$ . In the  $t$ -th phase, we define  $\Psi_i^{(t)} = \mathbf{1}(\|\Delta w_i^{(t)}\| > c) \cdot (\|\Delta w_i^{(t)}\| - c)$   
 251 as the incremental norm of the local update from  $f_i(w)$  compared to the clipping threshold  $c$ , for  
 252  $t = 1, 2, \dots, T$ .

253 In Definition [4.1](#), the incremental norm  $\Psi_i^{(t)}$  simply quantifies the difference between the norm of  
 254 the local update and its clipped version from  $f_i(w)$ . In the following, we will always assume the DP  
 255 noise injected  $\mathbb{E}[\|Q^{(t)}\|^2] = \sigma^2 d$ , following the classic privacy analysis of DP-SGD [\[38\]](#).

256 It is not hard to observe that the clipped local update is essentially a scaled version of the original  
 257 update, and thus virtually one may view DP-LSGD as a generalization of noisy LSGD but each local  
 258 update applies a different and adaptively-selected learning rate. To show meaningful characterization  
 259 on the difference among those learning rates, we need the following assumption as a generalization  
 260 of bounded-variance stochastic gradient.

261 **Assumption 4.1** (Incremental norm of Bounded Second Moment). *When applying the clipped version  
 262 of Algorithm 1 via 2 on an objective function  $F(w) = \frac{1}{n} \cdot f_i(w)$ ,  $\mathbb{E}[(\sum_{i=1}^n (\Psi_i^{(t)})^2)/n]$  is upper  
 263 bounded by  $\mathcal{B}^2$ , for some global parameter  $\mathcal{B}$  for  $t = 1, 2, \dots, T$ .*

264 Assumption 4.1 basically states that in expectation the square of  $l_2$ -norm of each local update is  
 265 bounded. Assumption 4.1 also suggests that  $\mathbb{E}[(\sum_{i=1}^n \Psi_i^{(t)})/n] \leq \mathcal{B}$ .

#### 266 4.1 Utility of DP-LSGD in Convex Optimization

267 Another assumption we need for the analysis of DP-LSGD on general convex optimization is the  
 268 similarity among the local functions.

269 **Assumption 4.2** ( $\gamma$  Similarity). *For  $F(w) = 1/n \cdot \sum_{i=1}^n f_i(w)$ , local functions  $f_i$  are of  $\gamma$ -similarity  
 270 to  $F$  such that for any  $w \in \mathcal{W}$ ,  $|f_i(w) - F(w)| \leq \gamma$ , for some constant  $\gamma > 0$ .*

271 The main reason why we need this additional Assumption 4.2 is because we do not assume Lipschitz  
 272 continuity of  $F(w)$ . Thus, we alternatively consider to use the similarity among local functions to  
 273 characterize the deviation of the evaluation of  $F(\cdot)$  on biased iterates.

274 **Theorem 4.1** (Last-iterate of DP-LSGD in General Convex Optimization). *For an arbitrary objective  
 275 function  $F(w) = \frac{1}{n} \cdot \sum_{i=1}^n f_i(w)$  where  $f_i(w)$  is convex and  $\beta$ -smooth, and under Assumptions 2.1  
 276 4.1 and 4.2 when  $\eta = O(1/\sqrt{TK})$  and  $Q^{(t)}$  is independent DP noise such that  $\mathbb{E}[Q^{(t)}] = 0$  and  
 277  $\mathbb{E}[\|Q^{(t)}\|^2] = \sigma^2 d$ ,  $t = 1, 2, \dots, T$ , then DP-LSGD with clipping threshold  $c$  ensures that*

$$\begin{aligned} \frac{c}{c + \mathcal{B}} \cdot \mathbb{E}[F(\bar{w}^{(T)}) - F(w^*)] &= \tilde{O}\left(\left(\frac{1}{\sqrt{TK}} + \frac{K}{nT}\right)\|\bar{w}^{(0)} - w^*\|^2\right. \\ &\quad \left.+ \left(\frac{K}{nT} + \frac{1}{\sqrt{TK}}\right)\left(1 + \frac{K^{3/2}}{\sqrt{T}} + \frac{K}{nq}\right)\tau + \left(\frac{K^{3/2}}{\sqrt{T}n} + 1\right)\frac{\gamma\mathcal{B}}{c + \mathcal{B}} + \sqrt{TK}\sigma^2 d\right). \end{aligned} \quad (5)$$

278 When  $K = O(nq)$  and  $K = O(T)$ , and for  $(\epsilon, \delta)$ -DP, where  $\sigma = \tilde{O}\left(\frac{c\sqrt{T\log(1/\delta)}}{n\epsilon}\right)$ , we have that

$$\begin{aligned} &\mathbb{E}[F(\bar{w}^{(T)}) - F(w^*)] \\ &= \tilde{O}\left(\underbrace{\frac{c + \mathcal{B}}{c} \cdot \left(\frac{\|\bar{w}^{(0)} - w^*\|^2}{\sqrt{TK}} + \left(\frac{1}{\sqrt{TK}} + \frac{K}{T}\right)\tau\right)}_{(A)} + \underbrace{\frac{\gamma\mathcal{B}}{c}}_{(B)} + \underbrace{\frac{c + \mathcal{B}}{c} \cdot \frac{T^{3/2}K^{1/2}\log(1/\delta)dc^2}{n^2\epsilon^2}}_{(C)}\right). \end{aligned}$$

279 The proof can be found in Appendix C. We focus on a practical scenario where  $\mathcal{B} = O(c)$ , i.e., the  
 280 incremental norm of local updates is in the same order of the clipping threshold  $c$  selected, and thus  
 281  $(c + \mathcal{B})/c = O(1)$ . From Theorem 4.1, we show the last-iterate utility of DP-LSGD is captured by  
 282 three terms: (A) a similar convergence rate as regular LSGD, (B) a clipping bias, and (C) the DP noise  
 283 variance. First, ignoring the bias and noise, DP-LSGD still enjoys a convergence rate  $\tilde{O}\left(\frac{\|\bar{w}^{(0)} - w^*\|^2}{\sqrt{TK}} + \left(\frac{1}{\sqrt{TK}} + \frac{K}{T}\right)\tau\right)$ , which is slightly worse compared to Theorem 3.2 with  $\tilde{O}\left(\frac{\|\bar{w}^{(0)} - w^*\|^2}{\sqrt{TK}} + \left(\frac{1}{\sqrt{TK}nq} + \frac{K}{T}\right)\tau\right)$  as a consequence of clipping which essentially applies different learning rates in each local  
 284 update. Second, the clipping bias is captured by  $(\gamma\mathcal{B})/c$ . This matches our intuition that a larger  
 285 incremental norm  $\mathcal{B}$  combined with a smaller clipping threshold  $c$  will imply a more significant change  
 286 on the local update and thus a larger bias. The last accumulated perturbation term is determined by  
 287 the noise injected across each phase with an effect of  $\tilde{O}\left(\frac{T^{3/2}K^{1/2}\log(1/\delta)dc^2}{n^2\epsilon^2}\right)$  for  $(\epsilon, \delta)$ -DP under  
 288  $T$ -fold composition.

292 As we consider the very generic setup with non-trivial clipping, Theorem 3.2 cannot be directly com-  
 293 pared to the classic DP-utility tradeoff [32] given Lipschitz continuity, where a utility loss  $\tilde{\Theta}(\sqrt{d}/n\epsilon)$

294 is tight for convex optimization under  $(\epsilon, \delta)$ -DP. However, we have the following interesting observa-  
 295 tions. First, when we take the clipping threshold  $c = O(\eta) = O(1/\sqrt{TK})$  and  $K = O(T \cdot d/(n^2\epsilon^2))$ ,  
 296 DP-LSGD achieves the same optimal rate  $\tilde{O}(\sqrt{d}/n\epsilon)$  [33] ignoring the clipping bias. Second and  
 297 more important, when the stochastic gradient variance  $\tau$  is in the same order of the clipping bias  
 298  $O(\gamma\mathcal{B}/c)$ , then by selecting  $c = \Theta(\eta)$  and  $K = \Theta(T)$ , Theorem 4.1 suggests that DP-LSGD will  
 299 converge in  $O(1/T)$  to an  $O(\gamma\mathcal{B}/c + \frac{d}{n^2\epsilon^2})$  neighborhood of the global optimum. As a comparison,  
 300 when we select  $K = 1$  in Theorem 4.1, it becomes the analysis of DP-SGD but the convergence  
 301 rate to the neighborhood of global optimum in the same scale  $O(\gamma\mathcal{B}/c + \frac{d}{n^2\epsilon^2})$  is only  $O(1/\sqrt{T})$ .  
 302 Moreover, as we will show in the next section, the local update bound  $\mathcal{B}$  in DP-SGD with  $K = 1$   
 303 in practice would be much larger than that of DP-LSGD with a relatively larger  $K$ . As a simple  
 304 generalization, we also include an analysis of DP-LSGD on strongly-convex functions in Appendix  
 305 D, and we move our focus to the non-convex optimization in the following.

## 306 4.2 Utility of DP-LSGD in Non-convex Optimization

307 **Theorem 4.2** (DP-LSGD in Non-convex Optimization). *For  $F(w) = \frac{1}{n} \cdot \sum_{i=1}^n f_i(w)$  where  $f_i(w)$   
 308 is  $\beta$ -smooth and satisfies Assumptions 2.1 and 4.1 when  $\eta = O(1/K)$ , DP-LSGD ensures that*

$$\mathbb{E}\left[\frac{\sum_{t=1}^T \|\nabla F(\bar{w}^{(t-1)})\|^2}{T}\right] \leq \frac{4F(\bar{w}^{(0)})}{TK\eta} + \frac{16\eta^2\tau\beta^2K^2}{nq} + \frac{4(1+\beta\eta)(\mathcal{B}^2/q + \sigma^2d)}{\eta^2K}. \quad (6)$$

309 When we select  $\eta = O(\frac{1}{\sqrt{TK}})$  and  $K = \Theta(T)$ , for  $(\epsilon, \delta)$ -DP we have that

$$\mathbb{E}\left[\frac{\sum_{t=1}^T \|\nabla F(\bar{w}^{(t-1)})\|^2}{T}\right] = \tilde{O}\left(\frac{F(\bar{w}^{(0)})}{T} + \frac{\tau}{nq} + \frac{\mathcal{B}^2T}{q} + \frac{d}{n^2\epsilon^2}\right). \quad (7)$$

310 The proof can be found in Appendix E. For the analysis of DP-LSGD in non-convex optimization,  
 we do *not* need Assumption 4.2 on the similarity among local functions and Theorem 4.2 is simply  
 obtained by substituting the clipping error from each phase into Theorem 3.2. To have a more clear  
 picture, we still consider a practical scenario when  $\mathcal{B} = \mathcal{B}_0 \cdot \eta$  for some constant  $\mathcal{B}_0$  and the variance  
 $\tau$  is also some constant. Then, from (7) we have that

$$\mathbb{E}\left[\frac{\sum_{t=1}^T \|\nabla F(\bar{w}^{(t-1)})\|^2}{T}\right] = O\left(\frac{F(\bar{w}^{(0)})}{T} + \frac{1}{nq} + \frac{\mathcal{B}_0^2}{q} + \frac{d}{n^2\epsilon^2}\right) = \tilde{O}\left(\frac{1}{T} + \frac{1}{q} + \frac{d}{n^2\epsilon^2}\right).$$

311 In other words, similar to the convex case, DP-LSGD will converge at a rate of  $O(1/T)$  to an  
 312  $\tilde{O}(1 + d/(n^2\epsilon^2))$  neighborhood of a saddle point given some constant sampling rate  $q$ . As a  
 313 comparison, for DP-SGD when  $K = 1$ , from Theorem 3.2 we can only ensure an  $O(1/\sqrt{T})$   
 314 convergence rate to a same  $\tilde{O}(1 + d/(n^2\epsilon^2))$  neighborhood.

## 315 5 Why DP-LSGD Produces Less Bias and Better SNR

316 Throughout the previous section, we showed that asymptotically DP-LSGD enjoys a faster conver-  
 317 gence rate to a neighborhood of (global/local) optimum compared to DP-SGD. We characterized  
 318 the clipping bias mainly based on the second moment upper bound  $\mathcal{B}^2$  of the incremental norm  
 319  $\Psi_i^{(t)}$  of local updates. In this section, we proceed to empirically study the  $\Psi_i^{(t)}$ , and the tradeoff  
 320 between clipping bias and DP (Gaussian) noise in practical deep learning tasks. We will explain why  
 321 DP-LSGD could produce smaller bias and enable more efficient clipping compared to DP-SGD.

322 To produce good utility-privacy tradeoff, a proper selection of the clipping threshold  $c$  is important.  
 323 Many existing works are devoted to optimizing the selection of  $c$  by either grid searching [35] or  
 324 adaptive fine-tuning [44]. A smaller  $c$  requires less DP noise. But, as a tradeoff shown in Theorem  
 325 4.1 and 4.2, a smaller  $c$  and a consequently a larger  $\mathcal{B}$  will also lead to a heavier clipping bias. Thus,  
 326 from the perspective of signal-to-noise ratio (SNR), an ideal scenario is that the  $l_2$ -norm of each  
 327 local update is *concentrated* such that we can maximize the efficiency of the clipping power  $c$  with  
 328 a small clipping effect for most local updates. Interpreted via our developed theory of clipping  
 329 bias, it is expected that given the clipping threshold  $c$ , the incremental norm  $\Psi_i^{(t)}$  would be small,  
 330 captured by  $\mathcal{B}$  in (5) and (7). In Fig. 1 (a,b), we plot various statistics of the incremental norm  $\Psi_i^{(t)}$

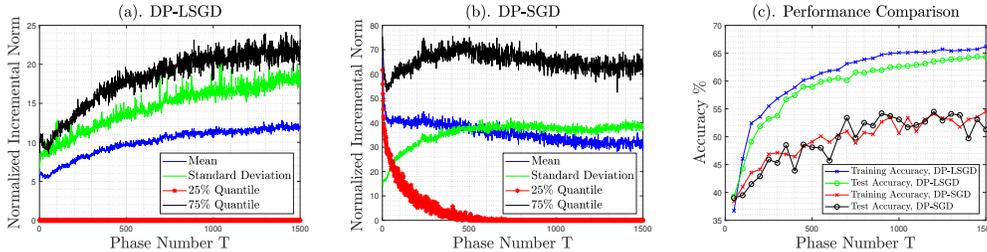


Figure 1: Training ResNet 20 on CIFAR10 with DP-LSGD ( $K = 10, \eta = 0.025, c = 1$ ) and DP-SGD ( $K = 1, \eta = 1, c = 1$ ) under  $(\epsilon = 2, \delta = 10^{-5})$ -DP, with expected batch size 1000.

331 for DP-LSGD and DP-SGD, respectively, on training CIFAR10 [45]. By our analysis, DP-LSGD  
 332 usually should apply a smaller learning rate  $\eta$ . To have a fair comparison, we consider the normalized  
 333 incremental norm  $\Psi_i^{(t)}/\eta$ . Given the same clipping threshold, comparing Fig. 1(a) and (b), the mean  
 334 of normalized incremental norm, captured by  $\mathcal{B}/\eta$  in our theorems, of DP-LSGD is only around 32%  
 335 compared to that of DP-SGD. The corresponding standard deviation is around only 40% compared to  
 336 that of DP-SGD. One may also compare the 25% and 75% quantiles, which suggest that more local  
 337 updates bear less clipping influence in DP-LSGD and thus enjoying a higher clipping efficiency. We  
 338 also report the comparison when training ResNet20 [46] on SVHN [47] in Fig. 2 in Appendix F with  
 339 similar observations. Details of experiment setups and the anonymous GitHub code link can be found  
 340 in Appendix F

Dataset and Method \ $\epsilon$	1.5	2.0	2.5	3.0	3.5	4.0
CIFAR10, DP-LSGD ( $K = 10$ )	59.4( $\pm 0.5$ )	64.0( $\pm 0.3$ )	66.2( $\pm 0.4$ )	67.7( $\pm 0.3$ )	68.7( $\pm 0.2$ )	69.9( $\pm 0.3$ )
CIFAR10, DP-SGD ( $K = 1$ )	49.8( $\pm 1.2$ )	58.7( $\pm 1.0$ )	59.9( $\pm 1.2$ )	60.6( $\pm 0.8$ )	62.1( $\pm 0.6$ )	62.8( $\pm 0.6$ )
SVHN, DP-LSGD ( $K = 10$ )	83.2( $\pm 0.4$ )	84.4( $\pm 0.5$ )	85.7( $\pm 0.5$ )	85.4( $\pm 0.4$ )	86.1( $\pm 0.4$ )	86.5( $\pm 0.3$ )
SVHN, DP-SGD ( $K = 1$ )	74.5( $\pm 0.8$ )	78.2( $\pm 0.6$ )	79.8( $\pm 0.6$ )	80.3( $\pm 1.0$ )	81.7( $\pm 0.4$ )	82.2( $\pm 0.5$ )

Table 1: **Test Accuracy** of ResNet20 on CIFAR10 and SVHN via DP-LSGD and DP-SGD under various  $\epsilon$  and fixed  $\delta = 10^{-5}$ , with expected batch size 1000.

341 In Fig. 1(c), we record the performance of DP-LSGD and DP-SGD, which coincides with our theory  
 342 that DP-LSGD has a smaller clipping bias and a faster convergence rate. The smaller incremental  
 343 norm in DP-LSGD is not surprising. With relatively larger  $K$ , for each individual function  $f_i(w)$ ,  
 344 though the  $K$  local gradients are correlated and essentially determined by a single sample, the  
 345 aggregation of them still averages out substantial sampling noise and makes the  $l_2$ -norm of local  
 346 updates more concentrated. In Table 1, we include additional comparison between their performance  
 347 on CIFAR10 [45] and SVHN [47]; DP-LSGD produces significant improvements.

## 348 6 Conclusion and Prospects

349 In this paper, via LSGD, we provide a unified analysis of the clipping bias and the utility loss in  
 350 privacy-preserving gradient methods for both centralized and distributed setups. Provided the generic  
 351 analysis, we develop the connections between the bias and the second moment of local updates.  
 352 This initializes a new direction to systematically instruct private learning by connecting the research  
 353 of variance reduction in distributed optimization. In this paper we only focus on regular LSGD  
 354 to show its advantage over DP-SGD, but advanced acceleration methods [30, 31, 43] are known  
 355 in non-private federated learning to further reduce the ‘‘local-update drift’’ caused by (per-sample)  
 356 data heterogeneity. This could then further reduce the clipping bias given local updates of smaller  
 357 variance. Thus, a promising future direction is to understand and incorporate those techniques  
 358 within the sensitivity control framework. Another important issue we have not fully explored is the  
 359 software implementation of DP-LSGD in the centralized case. For DP-SGD, many PyTorch libraries  
 360 with fast per-sample gradient computation in low memory overhead have been developed, such as  
 361 Opacus [48]. However, in all above-presented experiments, we simulate DP-LSGD in a distributed  
 362 environment and compute each local update in parallel at a cost of large memory. Given limited  
 363 hardware resources, this restricts the application of larger batchsize (tens of thousands) and deploying  
 364 deeper neural networks, which are known to produce much better utility-privacy tradeoffs [36, 49].  
 365 We leave empirical efficiency improvement to future work.

## References

- 366
- 367 [1] Sebastian Urban Stich. Local sgd converges fast and communicates little. In *ICLR 2019-*  
368 *International Conference on Learning Representations*, number CONF, 2019.
- 369 [2] Ahmed Khaled, Konstantin Mishchenko, and Peter Richtárik. Tighter theory for local sgd on  
370 identical and heterogeneous data. In *International Conference on Artificial Intelligence and*  
371 *Statistics*, pages 4519–4529. PMLR, 2020.
- 372 [3] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to  
373 sensitivity in private data analysis. In *Theory of Cryptography: Third Theory of Cryptography*  
374 *Conference, TCC 2006, New York, NY, USA, March 4-7, 2006. Proceedings 3*, pages 265–284.  
375 Springer, 2006.
- 376 [4] Graham Cormode, Somesh Jha, Tejas Kulkarni, Ninghui Li, Divesh Srivastava, and Tianhao  
377 Wang. Privacy at scale: Local differential privacy in practice. In *Proceedings of the 2018*  
378 *International Conference on Management of Data*, pages 1655–1658, 2018.
- 379 [5] Robin C Geyer, Tassilo Klein, and Moin Nabi. Differentially private federated learning: A  
380 client level perspective. *arXiv preprint arXiv:1712.07557*, 2017.
- 381 [6] Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh,  
382 and Dave Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv*  
383 *preprint arXiv:1610.05492*, 2016.
- 384 [7] Tao Lin, Sebastian U Stich, Kumar Kshitij Patel, and Martin Jaggi. Don’t use large mini-batches,  
385 use local sgd. In *International Conference on Learning Representations*, 2020.
- 386 [8] Cynthia Dwork. Differential privacy. In *Automata, Languages and Programming: 33rd*  
387 *International Colloquium, ICALP 2006, Venice, Italy, July 10-14, 2006, Proceedings, Part II 33*,  
388 pages 1–12. Springer, 2006.
- 389 [9] Jinshuo Dong, Aaron Roth, and Weijie J Su. Gaussian differential privacy. *Journal of the Royal*  
390 *Statistical Society Series B: Statistical Methodology*, 84(1):3–37, 2022.
- 391 [10] Cynthia Dwork, Aaron Roth, et al. The algorithmic foundations of differential privacy. *Founda-*  
392 *tions and Trends® in Theoretical Computer Science*, 9(3–4):211–407, 2014.
- 393 [11] Xiaokui Xiao and Yufei Tao. Output perturbation with query relaxation. *Proceedings of the*  
394 *VLDB Endowment*, 1(1):857–869, 2008.
- 395 [12] Cynthia Dwork, Guy N Rothblum, and Salil Vadhan. Boosting and differential privacy. In *2010*  
396 *IEEE 51st Annual Symposium on Foundations of Computer Science*, pages 51–60. IEEE, 2010.
- 397 [13] Shuang Song, Kamalika Chaudhuri, and Anand D Sarwate. Stochastic gradient descent with  
398 differentially private updates. In *2013 IEEE global conference on signal and information*  
399 *processing*, pages 245–248. IEEE, 2013.
- 400 [14] Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar,  
401 and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC*  
402 *conference on computer and communications security*, pages 308–318, 2016.
- 403 [15] Xiangyi Chen, Steven Z Wu, and Mingyi Hong. Understanding gradient clipping in private sgd:  
404 A geometric perspective. *Advances in Neural Information Processing Systems*, 33:13773–13782,  
405 2020.
- 406 [16] Xinwei Zhang, Xiangyi Chen, Mingyi Hong, Zhiwei Steven Wu, and Jinfeng Yi. Under-  
407 standing clipping for federated learning: Convergence and client-level differential privacy. In  
408 *International Conference on Machine Learning, ICML 2022*, 2022.
- 409 [17] Debraj Basu, Deepesh Data, Can Karakus, and Suhas Diggavi. Qsparse-local-sgd: Distributed  
410 sgd with quantization, sparsification and local computations. *Advances in Neural Information*  
411 *Processing Systems*, 32, 2019.

- 412 [18] Hao Yu, Rong Jin, and Sen Yang. On the linear speedup analysis of communication efficient  
413 momentum sgd for distributed non-convex optimization. In *International Conference on*  
414 *Machine Learning*, pages 7184–7193. PMLR, 2019.
- 415 [19] Jianyu Wang and Gauri Joshi. Cooperative sgd: A unified framework for the design and analysis  
416 of local-update sgd algorithms. *The Journal of Machine Learning Research*, 22(1):9709–9758,  
417 2021.
- 418 [20] Farzin Haddadpour and Mehrdad Mahdavi. On the convergence of local descent methods in  
419 federated learning. *arXiv preprint arXiv:1910.14425*, 2019.
- 420 [21] Blake Woodworth, Kumar Kshitij Patel, Sebastian Stich, Zhen Dai, Brian Bullins, Brendan  
421 McMahan, Ohad Shamir, and Nathan Srebro. Is local sgd better than minibatch sgd? In  
422 *International Conference on Machine Learning*, pages 10334–10343. PMLR, 2020.
- 423 [22] Huang Fang, Xiaoyun Li, Chenglin Fan, and Ping Li. Improved convergence of differential  
424 private sgd with gradient clipping. In *International Conference on Learning Representations*  
425 *2023*.
- 426 [23] Xiaodong Yang, Huishuai Zhang, Wei Chen, and Tie-Yan Liu. Normalized/clipped sgd with per-  
427 turbation for differentially private non-convex optimization. *arXiv preprint arXiv:2206.13033*,  
428 2022.
- 429 [24] LO Mangasarian. Parallel gradient distribution in unconstrained optimization. *SIAM Journal*  
430 *on Control and Optimization*, 33(6):1916–1925, 1995.
- 431 [25] Ryan McDonald, Keith Hall, and Gideon Mann. Distributed training strategies for the structured  
432 perceptron. In *Human language technologies: The 2010 annual conference of the North*  
433 *American chapter of the association for computational linguistics*, pages 456–464, 2010.
- 434 [26] Fan Zhou and Guojing Cong. On the convergence properties of a  $k$ -step averaging stochastic  
435 gradient descent algorithm for nonconvex optimization. *arXiv preprint arXiv:1708.01012*, 2017.
- 436 [27] Shiqiang Wang, Tiffany Tuor, Theodoros Salonidis, Kin K Leung, Christian Makaya, Ting  
437 He, and Kevin Chan. When edge meets learning: Adaptive control for resource-constrained  
438 distributed machine learning. In *IEEE INFOCOM 2018-IEEE conference on computer commu-*  
439 *nications*, pages 63–71. IEEE, 2018.
- 440 [28] Anastasia Koloskova, Nicolas Loizou, Sadra Boreiri, Martin Jaggi, and Sebastian Stich. A  
441 unified theory of decentralized sgd with changing topology and local updates. In *International*  
442 *Conference on Machine Learning*, pages 5381–5393. PMLR, 2020.
- 443 [29] Kevin Hsieh, Amar Phanishayee, Onur Mutlu, and Phillip Gibbons. The non-iid data quagmire  
444 of decentralized machine learning. In *International Conference on Machine Learning*, pages  
445 4387–4398. PMLR, 2020.
- 446 [30] Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and  
447 Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In  
448 *International Conference on Machine Learning*, pages 5132–5143. PMLR, 2020.
- 449 [31] Farzin Haddadpour, Mohammad Mahdi Kamani, Aryan Mokhtari, and Mehrdad Mahdavi.  
450 Federated learning with compression: Unified analysis and sharp guarantees. In *International*  
451 *Conference on Artificial Intelligence and Statistics*, pages 2350–2358. PMLR, 2021.
- 452 [32] Raef Bassily, Adam Smith, and Abhradeep Thakurta. Private empirical risk minimization:  
453 Efficient algorithms and tight error bounds. In *2014 IEEE 55th annual symposium on foundations*  
454 *of computer science*, pages 464–473. IEEE, 2014.
- 455 [33] Raef Bassily, Vitaly Feldman, Kunal Talwar, and Abhradeep Guha Thakurta. Private stochastic  
456 convex optimization with optimal rates. *Advances in neural information processing systems*, 32,  
457 2019.
- 458 [34] Nicolas Papernot, Abhradeep Thakurta, Shuang Song, Steve Chien, and Úlfar Erlingsson.  
459 Tempered sigmoid activations for deep learning with differential privacy. In *Proceedings of the*  
460 *AAAI Conference on Artificial Intelligence*, volume 35, pages 9312–9321, 2021.

- 461 [35] Florian Tramer and Dan Boneh. Differentially private learning needs better features (or much  
462 more data). In *International Conference on Learning Representations*, 2021.
- 463 [36] Soham De, Leonard Berrada, Jamie Hayes, Samuel L Smith, and Borja Balle. Unlock-  
464 ing high-accuracy differentially private image classification through scale. *arXiv preprint*  
465 *arXiv:2204.13650*, 2022.
- 466 [37] Yuqing Zhu and Yu-Xiang Wang. Poission subsampled rényi differential privacy. In *Internat-  
467 ional Conference on Machine Learning*, pages 7634–7642. PMLR, 2019.
- 468 [38] Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar,  
469 and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC*  
470 *conference on computer and communications security*, pages 308–318, 2016.
- 471 [39] Naman Agarwal, Ananda Theertha Suresh, Felix Xinnan X Yu, Sanjiv Kumar, and Brendan  
472 McMahan. cpsgd: Communication-efficient and differentially-private distributed sgd. *Advances*  
473 *in Neural Information Processing Systems*, 31, 2018.
- 474 [40] Tong Zhang. Solving large scale linear prediction problems using stochastic gradient descent  
475 algorithms. In *Proceedings of the twenty-first international conference on Machine learning*,  
476 page 116, 2004.
- 477 [41] Ohad Shamir and Tong Zhang. Stochastic gradient descent for non-smooth optimization:  
478 Convergence results and optimal averaging schemes. In *International conference on machine*  
479 *learning*, pages 71–79. PMLR, 2013.
- 480 [42] Xiaoyu Li and Francesco Orabona. On the convergence of stochastic gradient descent with  
481 adaptive stepsizes. In *The 22nd international conference on artificial intelligence and statistics*,  
482 pages 983–992. PMLR, 2019.
- 483 [43] Aritra Mitra, Rayana Jaafar, George J Pappas, and Hamed Hassani. Linear convergence in  
484 federated learning: Tackling client heterogeneity and sparse gradients. *Advances in Neural*  
485 *Information Processing Systems*, 34:14606–14619, 2021.
- 486 [44] Galen Andrew, Om Thakkar, Brendan McMahan, and Swaroop Ramaswamy. Differentially  
487 private learning with adaptive clipping. *Advances in Neural Information Processing Systems*,  
488 34:17455–17466, 2021.
- 489 [45] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.  
490 2009.
- 491 [46] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
492 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
493 pages 770–778, 2016.
- 494 [47] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng.  
495 Reading digits in natural images with unsupervised feature learning. 2011.
- 496 [48] Ashkan Yousefpour, Igor Shilov, Alexandre Sablayrolles, Davide Testuggine, Karthik Prasad,  
497 Mani Malek, John Nguyen, Sayan Ghosh, Akash Bharadwaj, Jessica Zhao, et al. Opacus:  
498 User-friendly differential privacy library in pytorch. *arXiv preprint arXiv:2109.12298*, 2021.
- 499 [49] Florian A Hölzl, Daniel Rueckert, and Georgios Kaissis. Equivariant differentially private deep  
500 learning. *arXiv preprint arXiv:2301.13104*, 2023.
- 501 [50] Moritz Hardt, Ben Recht, and Yoram Singer. Train faster, generalize better: Stability of  
502 stochastic gradient descent. In *International conference on machine learning*, pages 1225–1234.  
503 PMLR, 2016.