

# DIFFERENTIALLY PRIVATE CONSENSUS AND OPTIMIZATION ON COMMUNICATION CONSTRAINED DIRECTED GRAPHS

MONICA RIBERO, YIYUE CHEN AND HARIS VIKALO

**Abstract.** The rise of machine learning has been accompanied by growing privacy concerns and demands to protect users' sensitive data. At the same time, the amount of data that needs to be processed has been rapidly growing, bringing forth concerns related to limited communication resources available in practical settings. To this end, in this paper we study decentralized versions of consensus and convex optimization problems over directed graphs with communication and privacy constraints. Leveraging a local differential privacy model, we provide provable privacy guarantees for decentralized algorithmic frameworks that rely on sparsification to reduce the communication cost; while motivated by meeting communication constraints, sparsification is interpreted and exploited as a privacy amplification mechanism. To our knowledge, these are the first consensus and decentralized optimization frameworks that provide differential privacy for decentralized learning on directed graphs under communication constraints. The proposed scheme is tested on the consensus model with synthetic datasets, and a tag prediction model with logistic regression over a realistic Stackoverflow dataset. The experiments validate theoretical results and demonstrate efficacy of the proposed differentially private schemes.

**Key words.** Distributed optimization, Consensus, Machine Learning, Federated Learning, Communication efficiency, Differential Privacy

**AMS subject classifications.** 68Q25, 68R10, 68U05

**1. Introduction.** Decentralized consensus and convex optimization have been studied in a number of fields including machine learning, signal processing, and control [42, 47, 46]. They have emerged as attractive alternatives to centralized solutions limited by latency challenges [12, 26], high cost of communicating data to the central server [27], and, in many settings, privacy concerns that prohibit central data aggregation [24, 45, 48]. In consensus, a set of  $n$  nodes, each one with a data vector  $\mathbf{x}_i \in \mathbb{R}^d$ , for  $i \in [n] := \{1, \dots, n\}$ , aims to find the mean vector  $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$ . In decentralized optimization, the nodes collaborate to minimize the finite sum of local objective functions  $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$  over a convex compact constraint set  $\mathcal{X}$ , i.e., solve

$$(1.1) \quad \min_{\mathbf{x} \in \mathcal{X}} \left[ f(\mathbf{x}) := \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}) \right].$$

Recent multiagent applications over IoT networks, mobile devices and federated learning systems have reignited attention to problems wherein a number of nodes locally collects data and a peer-to-peer interaction allows them to estimate a parameter or optimize a function [8, 21]. However, to deliver potential benefits of decentralized solutions, two issues need to be addressed. First, even if data remains local, shared updates may be high dimensional, e.g., in federated learning where the updates are deep learning models and the number of parameters to be exchanged may be in the millions [18]; in such settings, limited energy and bandwidth typical of practical systems create a communication bottleneck. Second, in recent applications such as recommender systems, federated learning, and online learning, users' privacy may be compromised if sharing information potentially sufficient for identification [33, 1, 30]. Previous solutions to these problems have focused on fixed topologies as in federated learning where a server communicates with all agent nodes in the network – the setting equivalent to a fully connected star network topology [30]. More general decentralized optimization topologies that have previously been studied include undirected graphs

46 [46, 38, 52, 34, 14, 24, 28]. However, real-world networks are known to be time-  
 47 varying and support directed communication between the nodes; while there exist  
 48 communication-efficient protocols for decentralized consensus and optimization in such  
 49 setting [10], no prior work providing privacy guarantees therein exist.

50 In this work, we consider differential privacy [15, 16] – a statistical framework that  
 51 enables trade-off between data utility and privacy loss. The privacy loss of a query  
 52 to a database is defined as the probability of identifying an individual record in the  
 53 database from the output of the query. This requires trusting a central aggregator  
 54 to compute the query output and mechanism, which is not available in the fully  
 55 distributed and decentralized setting. We consider the local differential privacy model,  
 56 introduced first by [23], where each node has to protect its outputs by perturbing  
 57 any shared data or message. We propose (to our knowledge, the first) convergent  
 58 algorithms for decentralized consensus and optimization over directed graphs that  
 59 satisfy both communication and differential privacy constraints.

60 In particular, our contributions can be summarized as follows:

- 61 • We propose algorithms for decentralized consensus and convex optimization  
 62 over time-varying directed graphs with communication, and local and record-  
 63 level differential privacy constraints.
- 64 • We provide convergence analysis for both algorithms; for the consensus algo-  
 65 rithm, we establish linear convergence; for the optimization algorithm, we  
 66 show  $O(\ln T/\sqrt{T})$  convergence rate if the global objective function is convex  
 67 and  $O(\ln T/T)$  convergence rate if, in addition, local objective functions are  
 68 strongly convex.
- 69 • We provide a tight privacy analysis and show record-level differential pri-  
 70 vacy guarantees, with a utility-privacy trade-off of  $O\left(\frac{dn}{\epsilon r} + \frac{\sqrt{nd^3}}{\epsilon r}\right)$  for convex  
 71 functions, and  $O\left(\frac{p^2 nd^2}{\epsilon^2 r^2}\right)$  for strongly convex local objectives.
- 72 • We perform extensive numerical studies under various communication and  
 73 privacy settings, and investigate accuracy/communication/privacy trade-offs.

74 **Notation.** We represent vectors by lowercase bold letters and matrices by upper-  
 75 case letters.  $[A]_{ij}$  represents the  $(i, j)$  element of matrix  $A$ .  $\|\cdot\|$  represents the standard  
 76 Euclidean norm. For convenience, the symbols used in the paper are summarized in a  
 77 table in the supplementary document, Sec A.

78 **1.1. Related work and significance.** Prior work on consensus algorithms  
 79 considered both directed and undirected graphs [5, 44], as well as time-varying graphs  
 80 [20, 53, 43, 7]. By leveraging compressed communication, [24] developed the first  
 81 linearly convergent communication-efficient consensus algorithm over undirected time-  
 82 invariant graphs. [10] established the same type of results for directed time-varying  
 83 graphs. Prior work on decentralized optimization includes a gradient descent scheme  
 84 [38], the alternating direction method of multipliers (ADMM) [52], and decentralized  
 85 dual averaging methods [14, 34]. More recently, [28] studied how decentralization  
 86 and asynchronous SGD affect convergence. [24, 45] introduced a decentralized convex  
 87 scheme with limited communication and convergence guarantees; note that all of the  
 88 above prior work is limited to the undirected settings.

89 The subgradient-push [35] and Directed Distributed Gradient Descent (D-DGD)  
 90 [55] address optimization over directed graphs; they achieve  $O(\frac{\ln T}{\sqrt{T}})$  convergence rate.  
 91 In the same setting, [37] improve the convergence rate to linear for smooth and strongly  
 92 convex functions. However, these approaches do not consider limited communication  
 93 settings. [10] studies directed networks under communication constraints and proves

94  $O(\frac{\ln T}{\sqrt{T}})$  convergence rate. In the non-convex setting, [3] proposes Overlap Stochastic  
 95 Gradient Push which combines the push-sum algorithm with stochastic gradient  
 96 updates and proves the same sub-linear rate as in SGD.

97 An approach to reducing communication by taking advantage of local SGD  
 98 rounds is studied in [49]; this work is extended in [25] to the non-i.i.d. case. [50]  
 99 proposes MATCHA, an algorithm that improves over previous approaches and reduces  
 100 communication and computation costs by randomly sampling clients. They also show  
 101 that local updates in distributed optimization can accelerate convergence rate of the  
 102 algorithm.

103 Decentralized optimization under privacy constraints has been studied in [11],  
 104 which provides an overview of privatization mechanisms for data exchanged in net-  
 105 works. This works shows an impossibility of convergence, under a different privacy  
 106 model and where all messages in all iterations are noised. In contrast, we rely in the  
 107 post-processing property of differential privacy to design our algorithm and achieve a  
 108 better accuracy-privacy tradeoff. Authors in [19, 51] study consensus under different  
 109 definitions of privacy. The focus of our work is on providing differential privacy guar-  
 110 antees for decentralized consensus and optimization over communication-constrained  
 111 time-varying networks; related prior work includes [4] which incorporates differential  
 112 privacy guarantees but considers a setting wherein each user has a set of personalized  
 113 parameters. In contrast, the parameters in our problem are shared across nodes and  
 114 the aim is to achieve consensus or minimize the finite sum of local objective functions.

115 None of the above considers both privacy and communication constraints, often  
 116 present simultaneously in practice. *cpSGD*, introduced in [2], takes communication  
 117 and privacy into account by using efficient quantization via random rotations, and  
 118 introducing a binomial privacy mechanism that reduces the communication overhead.  
 119 That work is inspired by federated learning and the obtained guarantees are valid only  
 120 for the fully connected, undirected network topologies. Concurrent to our work, [9]  
 121 studied the trade-off between communication and privacy, but they only consider the  
 122 centralized model.

123 Finally, it is worth pointing out that satisfying differential privacy constraints  
 124 is typically more challenging in iterative settings, including those commonly used in  
 125 optimization algorithms; the iterative nature of such algorithms requires splitting  
 126 the privacy budget across iterations. [1, 54, 29] proposed privacy techniques that  
 127 account for noisy SGD which adds Gaussian noise to the gradient before updating  
 128 the parameters. This is further refined in [32] by leveraging Renyi differential privacy  
 129 [31], a relaxation of the traditional  $(\epsilon, \delta)$ -differential privacy. In our work we analyse  
 130 composition across iterations and dimensions using strong composition [22], that allows  
 131 for a more interpretable bound. In our experiments we use Renyi-DP to provide a  
 132 tighter, more realistic accounting of privacy.

133 *Organization.* We start by introducing some preliminary definitions in [subsec-](#)  
 134 [tion 2.1](#), followed by the consensus algorithm and analysis in [subsection 2.2](#). We  
 135 continue with optimization algorithms and analysis in [subsection 2.3](#), and experimental  
 136 results in [section 3](#). We include detailed proofs in [section 5](#) and finish with a discussion  
 137 in [section 6](#).

## 138 **2. Private and Communication Efficient Decentralized Algorithms.**

### 139 **2.1. Preliminaries.**

140 *Communication-constrained networks.* We model the connectivity in a network  
 141 with  $n$  nodes by a time-varying directed graph. At time  $t$ , the in-neighbor connectivity  
 142 matrix (row-stochastic),  $W_{in}^t$ , and the out-neighbor connectivity matrix (column-

143 stochastic),  $W_{out}^t$ , are defined as

$$144 \quad (2.1) \quad [W_{in}^t]_{ij} = \begin{cases} > 0, & j \in \mathcal{N}_{in,i}^t \\ 0, & \text{otherwise} \end{cases}, \quad [W_{out}^t]_{ij} = \begin{cases} > 0, & i \in \mathcal{N}_{out,j}^t \\ 0, & \text{otherwise} \end{cases}$$

145 where  $\mathcal{N}_{in,i}^t$  is the set of nodes that can send messages to node  $i$  (including  $i$ ) and  
 146  $\mathcal{N}_{out,j}^t$  is the set of nodes that can receive messages from node  $j$  (including  $j$ ) at time  
 147  $t$ . Node  $i$  knows both  $\mathcal{N}_{in,i}^t$  and  $\mathcal{N}_{out,i}^t$ , which is sufficient for the construction of  $W_{in}^t$   
 148 and  $W_{out}^t$ .

149 To comply with the communication constraints, the nodes in a network may need  
 150 to limit their communication to only a fraction of a full message, which we facilitate by  
 151 applying sparsification methods. To this end, let us introduce a sparsification operator  
 152  $Q : \mathcal{R}^d \rightarrow \mathcal{R}^d$ ; applying  $Q$  to a  $d$ -dimensional real vector returns a sparsified version of  
 153 that vector. If node  $i$  can communicate  $k$  out of  $d$  entries of a message, the probability  
 154 of any given entry actually being communicated is  $\frac{k}{d}$ .

155 *Differential privacy.* Differential privacy was first introduced in [15] as a mechanism  
 156 to prevent output queries to databases from disclosing the inclusion of a particular  
 157 record on the dataset. Formally, it is a bound on the probability of losing a record's  
 158 privacy by including it in the computation of a query.

DEFINITION 2.1 ((approximate) Differential Privacy). *We say that a randomized algorithm  $\mathcal{M}$  satisfies  $(\epsilon, \delta)$ -differential privacy  $((\epsilon, \delta)$ -DP) if for any pair of datasets  $\mathcal{D}$  and  $\mathcal{D}'$  differing by only one record and any subset of outcomes  $S \in \text{range}(\mathcal{M})$  it holds that*

$$Pr(\mathcal{M}(\mathcal{D}) \in S) \leq e^\epsilon \cdot Pr(\mathcal{M}(\mathcal{D}') \in S) + \delta.$$

159 THEOREM 2.2 (Theorem A.1 in [16]). *The  $\ell_2$ -sensitivity of a query  $f$  evaluated*  
 160 *on a dataset  $\mathcal{D}$  with range in  $\mathcal{R}^d$  is defined as  $\Delta_2(f) := \max_{\mathcal{D}, \mathcal{D}'} \|f(\mathcal{D}) - f(\mathcal{D}')\|_2$ ,*  
 161 *where  $\mathcal{D}$  and  $\mathcal{D}'$  are datasets differing in only one record. The Gaussian mechanism*  
 162 *with parameter  $\sigma$ ,  $G_{f,\sigma}$ , adds zero-mean Gaussian noise with variance  $\sigma^2$  to all  $d$*   
 163 *coordinates in query  $f$ ; formally,  $G_{f,\sigma}(\mathcal{D}) = f(\mathcal{D}) + N(0, \sigma^2 I_d)$ , allowing an abuse of*  
 164 *notation.  $G_{f,\sigma}(\mathcal{D})$  is  $(\epsilon, \delta)$ -DP if  $\sigma \geq \frac{\Delta_2(f)}{\epsilon} \sqrt{2 \log(\frac{1.25}{\delta})}$ .*

165 The Sampled Gaussian Mechanism with sampling rate  $p$  and noise variance  $\sigma^2$ ,  $SGM_{p,\sigma}$ ,  
 166 first generates a subset of  $\mathcal{D}$  by selecting points independently at random with proba-  
 167 bility  $p$ , computes the query on this random subset, and adds to it samples from a  
 168 zero-mean Gaussian distribution with variance  $\sigma^2$ . Further details and illustrations of  
 169 this mechanism are in the supplementary material, Sec B. Differential privacy assumes  
 170 there is a trusted central aggregator that possesses all users' data and computes private  
 171 output queries, allowing the aggregation to add less noise. In the decentralized setting  
 172 each user has its data locally, thus all users' outcomes have to be noised, making the  
 173 problem harder. In our case, all nodes share record-level differentially private messages.  
 174 This means that we protect users against an attack wishing to learn if a specific record  
 175 is in their data. To illustrate the different models, consider the setting where each  
 176 node  $i$  has  $r$  records and  $X_i$  is a query to node  $i$ , and  $\Delta_2(X_i) = \frac{1}{r}$ . Assume we want  
 177 to disclose the mean of  $X_1, X_2, \dots, X_n$ . In the central model, the sensitivity of query  
 178  $\bar{x}$  is  $\frac{1}{nr}$  so we only add gaussian noise with variance  $O(\frac{L}{\epsilon nr})$ . In the local model, we  
 179 would have to add for each node  $i$ , thus noise standard deviation is augmented by a  
 180 factor of  $n$ ,  $n = O(\frac{L}{\epsilon r})$ .

181 **2.2. Differentially-Private Communication-Efficient Consensus.** In gen-  
 182 eral, applying sparsification methods to existing consensus schemes, e.g. [6, 7, 35],

183 does not guarantee convergence since the sparsification operator causes non-vanishing  
 184 error. In [10], authors rely on entry-wise sparsification of a message vector and the  
 185 structure of the underlying connectivity matrices to propose a convergent consensus  
 186 algorithm. This is accomplished by splitting the vector-valued consensus problem into  $d$   
 187 scalar-valued sub-problems with connectivity matrices  $\{W_{in,m}\}_{m=1}^d$  and  $\{W_{out,m}\}_{m=1}^d$ ,  
 188 re-normalized according to the sparsification patterns. We assume node  $i$  has access  
 189 to vector  $\mathbf{x}_i$  and, following [10], introduce an auxiliary ‘‘surplus’’ vector  $\mathbf{y}_i \in \mathcal{R}^d$ ; at  
 190 iteration  $t$ ,  $\mathbf{y}_i^t$  records the change between consecutive state vectors,  $\mathbf{x}_i^t - \mathbf{x}_i^{t-1}$ . Both  $\mathbf{x}_i^t$   
 191 and  $\mathbf{y}_i^t$  can be communicated to the out-neighbors of node  $i$ . To simplify the notations,  
 192 we introduce  $\mathbf{z}_i^t \in \mathcal{R}^d$  defined as

$$193 \quad (2.2) \quad \mathbf{z}_i^t = \begin{cases} \mathbf{x}_i^t, & i \in \{1, \dots, n\} \\ \mathbf{y}_{i-n}^t, & i \in \{n+1, \dots, 2n\}. \end{cases}$$

194 Let  $Q(\mathbf{z}_i^t)$  denote a vector obtained by sparsifying  $\mathbf{z}_i^t$ , and  $[Q(\mathbf{z}_i^t)]_m$  be the  $m$ -th entry of  
 195  $Q(\mathbf{z}_i^t)$ . The re-normalization of the in-neighbor and out-neighbor connectivity matrices  
 196 is performed according to  
 (2.3)

$$197 \quad [A_m^t]_{ij} = \begin{cases} \frac{[W_{in}^t]_{ij}}{\sum_{j \in \mathcal{S}_m^t(i,j)} [W_{in}^t]_{ij}} & \text{if } j \in \mathcal{S}_m^t(i,j) \\ 0 & \text{otherwise,} \end{cases}, [B_m^t]_{ij} = \begin{cases} \frac{[W_{out}^t]_{ij}}{\sum_{i \in \mathcal{T}_m^t(i,j)} [W_{out}^t]_{ij}} & \text{if } i \in \mathcal{T}_m^t(i,j) \\ 0 & \text{otherwise,} \end{cases}$$

198 respectively, where  $\mathcal{S}_m^t(i,j) = \{j | j \in \mathcal{N}_{in,i}^t, [Q(\mathbf{z}_j^t)]_m \neq 0\} \cup \{i\}$  and  $\mathcal{T}_m^t(i,j) = \{i | i \in$   
 199  $\mathcal{N}_{out,j}^t, [Q(\mathbf{z}_i^t)]_m \neq 0\} \cup \{j\}$ . The connectivity and communication weights across the  
 200 network are summarized by mixing matrices; in particular, the  $m$ -th mixing matrix at  
 201 iteration  $t$  is defined as

$$202 \quad (2.4) \quad \bar{M}_m^t = \begin{bmatrix} A_m^t & \mathbf{0} \\ I - A_m^t & B_m^t \end{bmatrix},$$

203 Having defined  $\mathbf{z}_i^t$  and  $\bar{M}_m^t$ , respectively, we now introduce the update for  $\mathbf{z}_i^t$  in the  
 204 communication efficient and  $(\epsilon, \delta)$ -DP consensus algorithm (Algorithm 2.1):

$$205 \quad (2.5) \quad z_{im}^{t+1} = \sum_{j=1}^{2n} [\bar{M}_m^t]_{ij} [Q(\mathbf{z}_j^t)]_m + \mathbb{1}_{\{t \bmod \mathcal{B} = \mathcal{B} - 1\}} \gamma [F]_{ij} z_{jm}^{\mathcal{B}[t/\mathcal{B}]},$$

206 where  $F = \begin{bmatrix} \mathbf{0} & I \\ \mathbf{0} & -I \end{bmatrix}$  and  $m$  represents the coordinate index. Vectors  $\mathbf{z}_i^t$  are updated  
 207 according to the sparsification and multiplication of the mixing matrices at all time  
 208 steps except those that are multiples of  $\mathcal{B}$ , i.e.,

$$209 \quad (2.6) \quad t \bmod \mathcal{B} = \mathcal{B} - 1.$$

210 where  $\mathcal{B}$  is the window size parameter indicating that starting from any time  $t = k\mathcal{B}$   
 211 for all integers  $k \geq 0$ , the union graph over  $\mathcal{B}$  consecutive time steps forms a strongly  
 212 connected graph (See Assumption 2.4).

213 In the update (2.5), when the time step  $t$  satisfies (2.6), vectors  $\mathbf{z}_i^{\mathcal{B}[t/\mathcal{B}]}$ , stored  
 214 at time  $\mathcal{B}[t/\mathcal{B}]$ , are also included in the update. The term  $\sum_{j=1}^{2n} [F]_{ij} z_{jm}^{\mathcal{B}[t/\mathcal{B}]}$  in the  
 215 update facilitates to form the following update over  $\mathcal{B}$  time steps:

$$216 \quad (2.7) \quad \mathbf{z}_{im}^{(k+1)\mathcal{B}} = \sum_{j=1}^{2n} [\bar{M}_m((k+1)\mathcal{B} - 1 : k\mathcal{B}) + \gamma F]_{ij} Q(\mathbf{z}_j^{k\mathcal{B}})_m$$

217 , where  $\bar{M}_m((k+1)\mathcal{B}-1 : k\mathcal{B})$  represents the product of mixing matrices from time  
 218  $k\mathcal{B}$  to time  $(k+1)\mathcal{B}-1$ . Since neither any single mixing matrix nor this product  
 219 can ensure a non-zero spectral gap,  $\gamma F$  is added to ensure a non-zero spectral gap.  
 220 In addition,  $F$  has its  $(1, 1)$  and  $(2, 1)$  block equal to zero matrix and the rest two  
 221 diagonal blocks and therefore the stored vectors are not communicated, leading to the  
 222 update (2.5) above.

223 Our proposed  $(\epsilon, \delta)$  differentially private procedure is summarized as [Algorithm 2.1](#);  
 224 it builds upon [10] to incorporate the Gaussian mechanism and applies it with sparsifi-  
 225 cation. As we prove in [Theorem 2.6](#), adding privacy guarantees has no detrimental  
 226 effect on the convergence of decentralized consensus with sparsified updates, while  
 227 [Theorem 2.3](#) establishes that all messages guarantee  $(\epsilon, \delta)$  differential privacy.

228 **2.2.1. Privacy guarantees.** In the consensus case, we assume the ‘‘honest-but-  
 229 curious’’ security model [40] where all nodes execute the algorithm honestly but may  
 230 attempt to learn additional information from the received messages. In the first  
 231 iteration, each node receives a perturbed version of individual data vectors from its  
 232 neighbors.

233 Our adopted privacy mechanism adds a zero-mean Gaussian noise to query outputs,  
 234 i.e., user vectors, where the noise standard deviation is proportional to the  $L_2$ -sensitivity  
 235 of the query.

236 In the first iteration of our consensus algorithm, the sensitivity of the query at node  
 237  $i$  is  $\Delta_2(\mathbf{z}_i^1) = \sqrt{d}C$ , where  $C$  denotes a bound on the magnitude of the components of  
 238 the original message  $\mathbf{x}_i^0$ .

239 Note that in the consensus problem it is only necessary to add noise to the initial  
 240 vector, before the first iteration of the algorithm. This is the key insight to our  
 241 differentially private algorithm: posterior messages do not touch the data again.

242 This is because subsequent messages, being functions of the initial noisy message,  
 243 are already privatized by the post-processing property of differential privacy [16]. This  
 244 allows us to achieve a better accuracy-privacy trade-off than [39], [11].

---

**Algorithm 2.1** Communication Efficient and  $(\epsilon, \delta)$ -DP Consensus Algorithm

---

- 1: **Input:** Time horizon  $T$ , Initial state  $\mathbf{x}^0$ , Initialize  $\mathbf{y}^0 = \mathbf{0}$ , noise variance  $\sigma^2$ ,  
network connectivity parameter  $\mathcal{B}$  and  $\gamma$
  - 2: Noise initial data  $\mathbf{x}_i^0 \leftarrow \mathbf{x}_i^0 + \mathbf{b}_i$  for  $\mathbf{b}_i \sim N(0, \sigma^2 I_{d \times d})$ ,  $i = 1, \dots, n$
  - 3: Initialize  $\mathbf{z}^0$
  - 4: **for**  $t = 1, \dots, T$  **do**
  - 5:   Generate non-negative matrices  $W_{in}^t, W_{out}^t$
  - 6:   **for**  $m = 1, \dots, d$  **do**
  - 7:     construct a row-stochastic  $A_m^t$  and a column-stochastic  $B_m^t$  according to (2.3)
  - 8:     construct  $\bar{M}_m^t$  according to definition (2.4)
  - 9:     **for**  $i = 1, \dots, 2n$  **do**
  - 10:        $z_{im}^{t+1} = \sum_{j=1}^{2n} [\bar{M}_m^t]_{ij} [Q(\mathbf{z}_j^t)]_m + \mathbb{1}_{\{t \bmod \mathcal{B} = \mathcal{B}-1\}} \gamma [F]_{ij} z_{jm}^{\mathcal{B}\lfloor t/\mathcal{B} \rfloor}$
  - 11:     **end for**
  - 12:   **end for**
  - 13: **end for**
  - 14: **Result:** Local consensus values  $\mathbf{z}_i^T$  for nodes  $i = 1, \dots, n$
- 

245 Before starting the loop, we define  $z_j^0 = [\mathbf{x}_j^0 + \mathbf{b}_j, y_j^0]$ , where  $\mathbf{b}_j \sim N(0, \sigma^2)$  in the

246 first iteration of the algorithm, node  $i$  receives

$$247 \quad (2.8) \quad z_{im}^1 = \sum_{j=1}^{2n} [M_m^t]_{ij} [Q(z_j^t)]_m + N(0, \sigma^2).$$

248 To protect user  $j$  we only need to ensure sensitivity of  $z_{im}^1$ , the first aggregated message  
249 received by  $i$  from other nodes.

250 **THEOREM 2.3.** *Set  $\sigma = \frac{\sqrt{d}C}{\epsilon} \sqrt{2 \log(\frac{1.25}{\delta})}$ , Algorithm 2.1 is  $(\epsilon, \delta)$ - differentially*  
251 *private.*

252 *Proof.* Note that the sensitivity of  $\mathbf{x}_i^0$  is  $\Delta_2(\mathbf{x}_i^0) = \sqrt{d}C$ , thus (2.8) states the  
253 Gaussian mechanism for the first step and hence this iteration is  $(\epsilon, \delta)$ -DP. Since  
254 further iterations depend only on the original DP queries (without further access to  
255 raw data), the overall  $(\epsilon, \delta)$ -DP of the algorithm follows from the post-processing  
256 property of differential privacy [16].  $\square$

257 **2.2.2. Convergence guarantees.** Having provided privacy guarantees for Al-  
258 gorithm 2.1, we next study its convergence properties. In particular, we prove that the  
259 addition of privacy mechanism does not adversely affect convergence rate. We begin  
260 by making assumptions needed for the analysis.

261 *Assumption 2.4.* The graph induced by sparsification is  $\mathcal{B}$ -jointly connected, i.e.,  
262 starting from any time step  $t = k\mathcal{B}$  where  $k = 0, 1, \dots$ , the union graph over  $\mathcal{B}$   
263 consecutive time steps,  $\bigcup_{t=k\mathcal{B}}^{(k+1)\mathcal{B}-1} \mathcal{G}(t)$  is a strongly connected directed graph.

264 This is a common assumption for algorithms on directed networks [37].

265 *Assumption 2.5.* Given  $\gamma \in (0, 1)$ , the set of all possible mixing matrices  $\{M_m^t\}$ ,  
266  $\mathcal{U}_M$ , is finite.

267 We are now ready state the convergence result. We differ all proofs to section  
268 [section 5](#)

**THEOREM 2.6.** *Suppose Assumption 2.4 and 2.5 hold. Fix*

$$\gamma \in (0, \min_m \left\{ \frac{1}{(20 + 8n)^n} (1 - |\lambda_3(\bar{M}_m((k+1)\mathcal{B} - 1 : k\mathcal{B}))|)^n \right\}),$$

269 *and let  $\tau = \max_{C \in \mathcal{U}_M} |\lambda_2(C)| < 1$ ,  $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^0$ , and  $t \geq 0$ . Then running*  
270 *Algorithm 2.1 for  $t$  iterations, suppose  $t = k\mathcal{B} - 1 + t'$ , where  $t' = 0, \dots, \mathcal{B} - 1$  and it*  
271 *holds that for any  $i \in [n]$  and  $t \geq 1$ ,*

$$272 \quad (2.9) \quad \begin{aligned} \|\mathbf{x}_i^t - \bar{\mathbf{z}}^t\| &\leq \sqrt{2nd} (\tau^{1/\mathcal{B}})^{t-(t'-1)} \sum_{j=1}^{2n} \sum_{m=1}^d |z_{jm}^0|, \\ \|\mathbf{y}_i^t\| &\leq \sqrt{2nd} (\tau^{1/\mathcal{B}})^{t-(t'-1)} \sum_{j=1}^{2n} \sum_{m=1}^d |z_{jm}^0|, \end{aligned}$$

273 *where  $\bar{\mathbf{z}}^t = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^t + \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i^t$ . Further,  $E[\mathbf{x}_i^t]$  converges to  $\bar{\mathbf{x}}$  at a linear rate*  
274  *$O(\tau^{t/\mathcal{B}})$ .*

275 **Theorem 2.6** implies that the local state vector  $\mathbf{x}_i^t$  converges (at a linear rate)  
276 to the averaging consensus vector  $\bar{\mathbf{z}}^t$ , while the surplus vectors  $\mathbf{y}_i^t$  vanishes to zero.

277 Since the noise added in the first iteration is unbiased, it always holds that  $E[\bar{\mathbf{z}}^t] = \bar{\mathbf{z}}^0$ .  
 278 Therefore, by setting  $\mathbf{y}_i^0 = \mathbf{0}$  in the initialization of Algorithm 2.1,  $\bar{\mathbf{z}}^0 = \bar{\mathbf{x}}$  and we  
 279 are able to guarantee linear convergence to  $\bar{\mathbf{x}}$  in expectation. Note that the fastest  
 280 consensus algorithms for directed time-varying graphs, be they for full communication  
 281 as in the [13, 6] or with sparsification of messages as in [10], enjoy linear convergence.  
 282 Unlike our Algorithm 2.1, however, those schemes do not provide privacy guarantees.

283 **2.3. Differentially-Private Communication-Efficient Optimization.** Un-  
 284 like consensus, decentralized optimization typically requires algorithms to access local  
 285 raw data whenever gradients are computed; consequently, all iterations with gradient  
 286 computation need to be protected. To this end, as in DP-SGD we deploy a privacy  
 287 mechanism by perturbing the gradient

$$(2.10) \quad \mathbf{g}_i^t = \begin{cases} \nabla f_i(\mathbf{x}_i^t), & i \in \{1, \dots, n\} \\ \mathbf{0}, & i \in \{n+1, \dots, 2n\} \end{cases}$$

289 with a noise term sampled from a Gaussian  $N(0, \sigma^2 D^2)$  distribution where  $D$  is an  
 290 upper bound on the magnitude of the components of  $\mathbf{g}_i^t$ . Note that (2.10) implies the  
 291 state vectors are updated via decentralized gradient descent while the surplus vectors  
 292 are updated the same way as in Algorithm 2.1.

293 Our proposed differentially private decentralized optimization scheme is formalized  
 294 as Algorithm 2.2. Differential privacy of the iterates is enforced in line 8 of the  
 295 algorithm. Despite incorporating privacy mechanism by adding noise to gradients,  
 296 the algorithm converges when a decreasing learning rate is used. Formal privacy and  
 297 convergence guarantees are provided below. Further, we study the trade-off between  
 298 these two.

299 **2.3.1. Privacy guarantees.** While the consensus algorithm only access the  
 300 original data on the first iteration, gradient descent access it at every iteration to  
 301 compute the current gradient. We have to take again a local DP approach but now  
 302 ensure that each message is privatized, and use the composition theorem to ensure the  
 303 track the overall algorithm DP guarantees.

304 Traditionally, the full gradient is considered as one query and privatized with a  
 305 gaussian vector  $b_t$  drawn from  $b_t \sim \mathcal{N}(0, \sigma^2 I_d)$ . Notice the expected norm of the noise  
 306 is  $\mathbb{E}[\|b_t\|_2] = \sqrt{d}\sigma$ . Thanks to the sparsification, we only have to take into account  
 307  $(1-q)d$  dimensions, and get a privacy amplification factor of  $(1-q)$  and the privacy  
 308 guarantee.

309 **THEOREM 2.7.** *Assuming  $\sigma = O\left(\frac{D\sqrt{T(1-q)d\log(1/\delta)}}{\epsilon}\right)$ , after  $T$  iterations, Algo-*  
 310 *rithm 2.2 satisfies  $(\epsilon, \delta)$  differential privacy.*

311 *Proof.* Assuming  $|g_{im}^t| \leq D$ , which is readily achieved by the Lipschitz assumption,  
 312 the sensitivity of  $z_{im}^{t+1}$  is given by  $\Delta_2(\alpha_t g_{im}) = \alpha_t D$ . We use advanced composition  
 313 (See Corollary 1 in [31]) over dimensions  $(1-q)d$ , and over  $T$  iterations, and obtain  
 314 the desired result.  $\square$

315 **2.3.2. Convergence guarantees.** We now shift our attention from the privacy  
 316 to convergence guarantees. To facilitate the analysis, we first study a broader problem:  
 317 we focus on convergence under arbitrary gaussian noise variance  $\sigma$ . Then, we study  
 318 the utility - privacy trade-off by imposing necessary assumptions for privacy on the  
 319 noise variance, that guarantee we meet privacy requirements.

320 We can now state the convergence result.

---

**Algorithm 2.2** Communication-Efficient and  $(\epsilon, \delta)$ -DP Decentralized Optimization Algorithm

---

**Input:** Time horizon  $T$  Initial state  $\mathbf{x}^0$ , Initialize  $\mathbf{y}^0 = \mathbf{0}$ , sparsification level  $q$ , noise variance  $\sigma^2$ , network connectivity parameter  $\mathcal{B}$  and  $\gamma$   
 set  $\mathbf{z}^0 = [\mathbf{x}^0; \mathbf{y}^0]$   
**for**  $t = 1, \dots, T$  **do**  
   Generate non-negative matrices  $W_{in}^t, W_{out}^t$   
   **for**  $m = 1, \dots, d$  **do**  
     construct a row-stochastic  $A_m^t$  and a column-stochastic  $B_m^t$  according to 2.3  
     construct  $\bar{M}_m^t$  according to definition  
     **for**  $i = 1, \dots, 2n$  **do**

$$z_{im}^{t+1} = \sum_{j=1}^{2n} [\bar{M}_m^t]_{ij} [Q(\mathbf{z}_j^t)]_m + \mathbb{1}_{\{t \bmod \mathcal{B} = \mathcal{B}-1\}} \gamma [F]_{ij} z_{jm}^{\mathcal{B}\lfloor t/\mathcal{B} \rfloor} - \mathbb{1}_{\{t \bmod \mathcal{B} = \mathcal{B}-1\}} \alpha_{\lfloor t/\mathcal{B} \rfloor} (g_{im}^{\mathcal{B}\lfloor t/\mathcal{B} \rfloor} + N(0, \sigma^2)),$$

**end for**  
   **end for**  
**end for**

**Result:** Local optimum values  $\mathbf{z}_i^T$  for nodes  $i = 1, \dots, n$

---

321 THEOREM 2.8. Suppose Assumptions 2.4-2.5 on mixing matrices hold, and fix  
 322  $\gamma \in (0, \min_m \left\{ \frac{1}{(20+8n)^n} (1 - |\lambda_3(\bar{M}_m((k+1)\mathcal{B} - 1 : k\mathcal{B}))|)^n \right\})$ . Assume  $|g_{im}^t| \leq D$ . If  
 323 the objective function  $f$  is convex, and the step size  $\alpha_t$  decays as  $O(1/\sqrt{T})$ , Algorithm  
 324 2.2 converges at a rate of  $O(\ln T/\sqrt{T})$ , Concretely,

$$325 \quad (2.11) \quad (E[f_{\min, T}] - f^*) \leq \frac{C_1}{\sum_{k=0}^{\lfloor T/\mathcal{B} \rfloor} \alpha_k} + \frac{C_2 \sum_{k=0}^{\lfloor T/\mathcal{B} \rfloor} \alpha_k^2}{\sum_{k=0}^{\lfloor T/\mathcal{B} \rfloor} \alpha_k} \leq \frac{C_1}{\sqrt{T/\mathcal{B} - 1}} + \frac{C_2 \ln(T/\mathcal{B})}{\sqrt{T/\mathcal{B} - 1}}$$

326 where  $f_{\min, T} := \min_{t=1, \dots, T} f(\bar{\mathbf{z}}^t)$  and  $f^*$  is the optimal value, and

$$327 \quad (2.12) \quad C_1 = \frac{n\sqrt{d}D^2}{2} + \frac{\sqrt{2}dD}{\sqrt{n}} \sum_{j=1}^{2n} \frac{\|\mathbf{z}_j^0\|}{1 - \tau^2},$$

328

$$329 \quad (2.13) \quad C_2 = \frac{(d + \sigma^2)nD^2}{2} + \frac{\sqrt{2}dD}{\sqrt{n}} \sum_{j=1}^{2n} \|\mathbf{z}_j^0\| + \frac{2\sqrt{2nd^2(d + \sigma^2)}D^2}{1 - \tau} + \frac{4\sqrt{d(d + \sigma^2)}D^2}{n}$$

330 THEOREM 2.9. Suppose Assumptions 2.4-2.5 on mixing matrices hold, fix  $\gamma \in$   
 331  $(0, \min_m \left\{ \frac{1}{(20+8n)^n} (1 - |\lambda_3(\bar{M}_m((k+1)\mathcal{B} - 1 : k\mathcal{B}))|)^n \right\})$ . Assume  $|g_{im}^t| \leq D$  and let  
 332  $\hat{\mathbf{x}}_i^T = \sum_{t=1}^{\lfloor T/\mathcal{B} \rfloor} \frac{(t-1)\mathbf{x}_i^t}{t(t-1)/2}$  for  $T \geq \mathcal{B}$  and  $K = \lfloor T/\mathcal{B} \rfloor$ . If the objective function  $f$  is convex  
 333 and the local objective function  $f_i$  is  $\mu_i$  strongly-convex, then there exists some constant

334  $C_3 > 0, C_4 > 0$  such that for all  $i$   
 (2.14)

$$E[f(\hat{\mathbf{x}}_i^T) - f(\mathbf{x}^*)] \leq \frac{C_3}{K} \sum_{j=1}^n \|\mathbf{x}_j^0\|_1 + \frac{C_4}{K} (1 + \ln(K-1)) + \frac{p^2 n D^2 d}{K} (1 + \sigma^2)$$

335

$$E\left[\sum_{j=1}^n \mu_j \|\hat{\mathbf{x}}_j^T - \mathbf{x}^*\|^2\right] \leq \frac{C_3}{K} \sum_{j=1}^n \|\mathbf{x}_j^0\|_1 + \frac{C_4}{K} (1 + \ln(K-1)) + \frac{p^2 n D^2 d}{K} (1 + \sigma^2)$$

336 where the step size  $\alpha_t = \frac{p}{t}$  for  $p \geq \frac{4n}{\sum_{i=1}^n \mu_i}$ ,  $\mathbf{x}^*$  is the optimal solution,  $C_3 = \frac{5\sqrt{2}nndD}{1-\tau}$

337 and  $C_4 = \frac{5\sqrt{2}nd(1+\sigma^2)n^2 dD^2}{1-\tau}$ .

338 **THEOREM 2.10.** Suppose Assumptions 2.4-2.5 on mixing matrices hold for a specific  $\gamma$ . Assume  $|g_{im}^t| \leq D = O(1)$ . If the objective function  $f$  is convex, and the  
 339 step size  $\alpha_t$  decays as  $O(1/\sqrt{T})$ , Algorithm 2.2 converges at a rate of  $O(\ln T/\sqrt{T})$ .  
 340 Concretely,  
 341

$$342 \quad (2.15) \quad (E[f_{\min,T}] - f^*) \leq O\left(\frac{C_1}{\sqrt{T}} + \frac{C_2 \ln T}{\sqrt{T}}\right)$$

343 where  $f_{\min,T} := \min_{t=1,\dots,T} f(\bar{\mathbf{z}}^t)$  and  $f^*$  is the optimal value, and

$$344 \quad (2.16) \quad C_1 = O(n\sqrt{d} + \sqrt{nd}), \quad C_2 = O\left((d + \sigma^2)n + \sqrt{nd}D + \sqrt{nd^2(d + \sigma^2)}\right)$$

345 **THEOREM 2.11.** Suppose Assumptions 2.4-2.5 on mixing matrices hold, fix  $\gamma \in$   
 346  $(0, \min_m \left\{ \frac{1}{(20+8n)^n} (1 - |\lambda_3(\bar{M}_m((k+1)\mathcal{B} - 1 : k\mathcal{B})|)|)^n \right\})$ . Assume  $|g_{im}^t| \leq D$  and let  
 347  $\hat{\mathbf{x}}_i^T = \frac{\sum_{t=1}^{\lfloor T/\mathcal{B} \rfloor} (t-1)\mathbf{x}_i^t}{t(t-1)/2}$  for  $T \geq \mathcal{B}$  and  $K = \lfloor T/\mathcal{B} \rfloor$ . If the objective function  $f$  is convex  
 348 and the local objective function  $f_i$  is  $\mu_i$  strongly-convex, then there exists some constant  
 349  $C_3 > 0, C_4 > 0$  such that for all  $i$   
 (2.17)

$$E[f(\hat{\mathbf{x}}_i^T) - f(\mathbf{x}^*)] \leq \frac{C_3}{K} \sum_{j=1}^n \|\mathbf{x}_j^0\|_1 + \frac{C_4}{K} (1 + \ln(K-1)) + \frac{p^2 n D^2 d}{K} (1 + \sigma^2)$$

350

$$E\left[\sum_{j=1}^n \mu_j \|\hat{\mathbf{x}}_j^T - \mathbf{x}^*\|^2\right] \leq \frac{C_3}{K} \sum_{j=1}^n \|\mathbf{x}_j^0\|_1 + \frac{C_4}{K} (1 + \ln(K-1)) + \frac{p^2 n D^2 d}{K} (1 + \sigma^2)$$

351 where the step size  $\alpha_t = \frac{p}{t}$  for  $p \geq \frac{4n}{\sum_{i=1}^n \mu_i}$ ,  $\mathbf{x}^*$  is the optimal solution,  $C_3 = \frac{5\sqrt{2}nndD}{1-\tau}$

352 and  $C_4 = \frac{5\sqrt{2}nd(1+\sigma^2)n^2 dD^2}{1-\tau}$ .

353 **Remark 2.12.** It is of interest to explore the impact of  $q$  (defined in) on the  
 354 convergence speed. As the mixing matrices are constructed over sparsified graphs,  $q$   
 355 affects the number of non-zero entries in the matrix and further affects the second  
 356 largest magnitude of eigenvalues. Specifically, when the graph connectivity parameter  
 357  $\mathcal{B}$  is fixed, greater  $q$  leads to greater  $\tau$  and further slows down the convergence process.

358 Theorem 2.8 and 2.9 provide convergence results for Algorithm 2.2 with different  
 359 assumptions: convexity of the global function (Theorem 2.8) and, in addition, strong-  
 360 convexity of local functions (Theorem 2.9). The resulting convergence rates match  
 361 those of the full communication gradient-push and D-DGD algorithms [35, 55], the

362 communication-efficient algorithm in [10], and the stochastic gradient-push [36], under  
 363 respective assumptions.

364 In the following section we study relevant instances where these assumptions hold,  
 365 including linear regression and logistic regression.

366 **2.3.3. Utility - privacy tradeoff.** In this section we provide explicit trade-offs  
 367 between privacy and utility of optimization algorithms. We state our results in ??

368 COROLLARY 2.13. *Assume the setting of Theorem 2.8 holds, particularly,  $f$  is a*  
 369  *$D$ -Lipschitz function, and assume  $D = \|\mathbf{z}_j^0\| = O(1)$ . Let  $r$  be the minimum number*  
 370 *of records each node has. Setting  $\sigma = O\left(\frac{\sqrt{T(1-q)d \log(1/\delta)}}{\epsilon r}\right)$ , after  $T = \epsilon^2 r^2$  iterations*  
 371 *algorithm Algorithm 2.2 is  $(\epsilon, \delta)$ - differentially private and the empirical risk is bounded*  
 372 *by*

$$373 \quad E[f_T - f^*] \leq O\left(\frac{dn}{\epsilon r} + \frac{\sqrt{nd^3}}{\epsilon r}\right)$$

374 COROLLARY 2.14. *Assume the same setting of Theorem 2.9, and  $r$  be the minimum*  
 375 *number of records each user has. Let  $\sigma = O\left(\frac{\sqrt{T(1-q)d \log(1/\delta)}}{\epsilon r}\right)$ , then Algorithm 2.2*  
 376 *as  $T \rightarrow \infty$ ,*

$$377 \quad E[f_T - f^*] \leq O\left(\frac{p^2 nd^2 (1-q)}{\epsilon^2 r^2}\right)$$

378 The proof of both corollaries follows by replacing  $\sigma$  in Theorem 2.8 and ?? with  
 379 the appropriate value of  $\sigma$ .

380 **3. Numerical Results.** In this section, we demonstrate performance of the  
 381 proposed privacy-preserving algorithms for decentralized consensus and optimization.  
 382 In both settings we show that, as expected, privacy and communication constraints slow  
 383 down convergence but the developed methods ultimately achieve performance similar  
 384 to that of non-private and full-communication algorithms. We start our numerical  
 385 studies with a network system having 10 nodes, and generate its edges randomly while  
 386 preserving the strong connectivity.

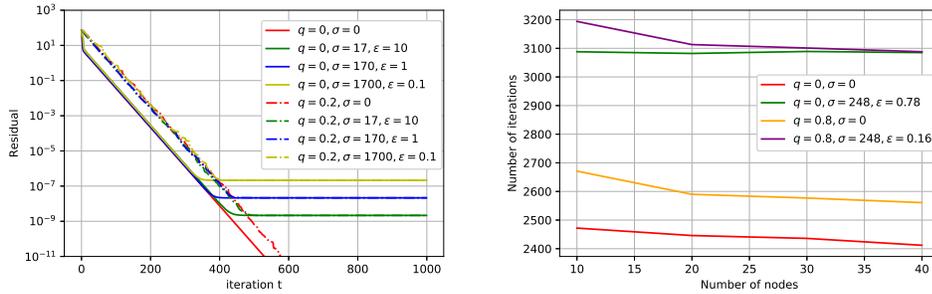
387 The construction begins with the Erdős–Rényi model [17] with edge probability  
 388 parameter equal to 0.9; then, 2 directed edges are dropped from each strongly con-  
 389 nected graph, leading to directed graphs. Building upon this basic structure, we can  
 390 construct networks with different connectivity. Recall that the window size parameter  
 391  $\mathcal{B}$ , introduced in Assumption 2.4, implies that the union graph over  $\mathcal{B}$  consecutive time  
 392 steps, starting from any time that is a multiple of  $\mathcal{B}$ , forms an almost-surely strongly  
 393 connected Erdős–Rényi graph. When  $\mathcal{B} = 1$ , the network is strongly connected at each  
 394 time step. We then apply sparsification such that the communication throughput is  
 395 brought down to various sparsity levels  $q$  (larger  $q$  means more entries are sparsified,  
 396  $q = 0$  means full communication). For privacy accounting in optimization we use the  
 397 TensorFlow Privacy library.<sup>1</sup>

398 **3.1. Consensus.** In the consensus problem, each node has access to a local  
 399 vector of dimension  $d = 64$ . Components of the initial local vector at node  $i$ ,  $\mathbf{x}_i^0$ , are  
 400 generated uniformly at random from  $[-5, 5]$ . To illustrate the effect of the privacy  
 401 mechanism, in Figure 1 we compare the performance of our Algorithm 1 for different  
 402 levels of noise  $\sigma$  and sparsity  $q$ , and show the corresponding privacy guarantee  $\epsilon$ . In  
 403 Fig. 1a, we show the residual as a function of the number of iterations  $t$ . We observe  
 404 that sparsity and noise added to provide privacy only delay the convergence without

<sup>1</sup><https://github.com/tensorflow/privacy>

405 affecting its rate, matching the results of Theorem 2.6. As expected, higher values of  
 406  $q$  result in slower convergence since less information is communicated; however, higher  
 407  $q$  achieves higher privacy for fixed  $\sigma$  because the probability of observing an entry  
 408 is lower. Fig. 1b shows that for a fixed sparsification level, the convergence becomes  
 409 faster as the number of nodes increases. Further results for varied values of parameters  
 410 and network topologies are in the supplementary material, Sec C.

411 We observe the noise effectively does not affect the final residual, neither the  
 412 rate of convergence: in Figure 1 it is clear that for all values of  $q$  and  $\sigma$  the speed of  
 413 convergence is the same, although the initial residual might differ depending on the  
 414 communication and privacy parameters. Eventually all methods converge.



(a) Residual vs. iterations for a 10-node network,  $\mathcal{B} = 5$  and  $\delta = 10^{-4}$ . (b) The number of iterations needed for residual to drop below  $10^{-10}$  as a function of the number of nodes.

Fig. 1: Convergence of Algorithm 1 for varied parameters  $q$  and  $\sigma$ , and the privacy loss  $\epsilon$  achieved. In (a), we see sparsity and noise delay convergence in early iterations but the convergence rate is unaffected. In (b) we show that for a fixed sparsification level, the convergence becomes faster as the number of nodes increases.

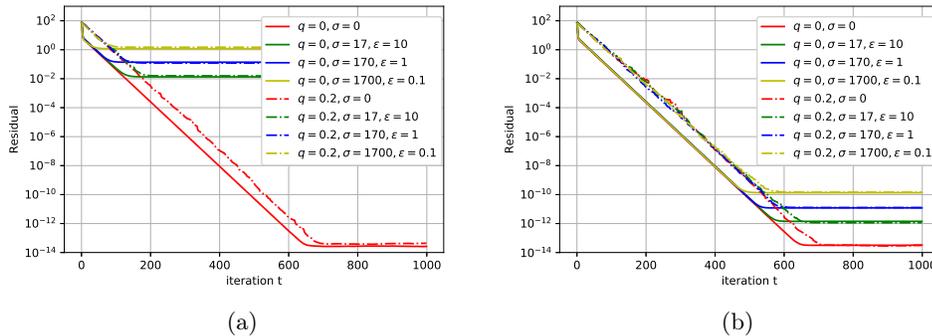


Fig. 2: Residual vs. iterations for a 10-node network,  $\mathcal{B} = 5$  and  $\delta = 10^{-4}$ .

415 **3.2. Decentralized Optimization Problems.** Next, we test performance of  
 416 Algorithm 2 on a multi-class tag classification task with a logistic regression model,  
 417 which leads to the optimization problem with features  $\mathbf{m}_{ij}$  and corresponding label  
 418  $y_{ij}$  of the form

$$419 \quad (3.1) \quad \min_{\mathbf{x}} \left\{ \frac{\mu}{2} \|\mathbf{x}\|^2 + \sum_{i=1}^n \sum_{j=1}^N \ln(1 + \exp(-(\mathbf{m}_{ij}^T \mathbf{x}_i) y_{ij})) \right\}.$$

420 with regularization parameter  $\mu$ .

421 The model is trained and tested on the Stackoverflow dataset, a language modelling  
 422 dataset with questions and answers collected from 342477 unique users. The objective  
 423 is to tag each sentence with appropriate categories. We present detailed preprocessing  
 424 of the dataset in the supplementary, Sec C. Following prior work [41], we use a build  
 425 vocabulary with 10000 frequent words and restrict each user’s dataset to have at most  
 426 128 sentences. We rely on padding and truncation to enforce 20 word sentences, and  
 427 represent them with index sequences corresponding to the vocabulary words, out of  
 428 vocabulary words, beginning and end of sentences.

429 The 150,000 data points are randomly split into 10 groups of equal size, where  
 430 each group is interpreted as being the local data for one of the nodes in the network.  
 431 Each node uses 13,500 data points as training data and the remaining 1500 points as  
 432 the validation set. The testing data contains 37640 data points.

433 We consider a network with 10 nodes and evenly split 150,000 data points at  
 434 random into 10 groups, each one representing one node in the network. We leave  
 435 1500 points for validation for each node. For this problem we use a noise variance of  
 436  $\sigma D = 30$ , a the step size  $\alpha_t = \frac{0.02}{t}$  and the privacy parameter  $\delta = 10^{-4}$ .

437 In ?? we observe the effect of privacy and sparsity. Both constraints slightly delay  
 438 and affect convergence. However, Figure 3a shows that this has minimal impact on  
 439 the accuracy, and that after a few rounds all models reach a similar level of accuracy.  
 440 Finally, Figure 3b shows that we are able to maintain a fair privacy budget ( $\epsilon < 10$ )  
 441 for models, even at the end of the training; this is a reasonable budget for iterative  
 442 procedures in literature [1], showing that our algorithms are able to achieve very good  
 443 performance while guaranteeing privacy and meeting communication constraints.

#### 444 4. Additional experiments.

445 **4.1. Consensus.** In Figure 1 of the main paper, we show convergence results for  
 446 the proposed consensus algorithm (algorithm 1). With the same parameter setups,  
 447 Figure 4 in this document illustrates the relationship between sparsity level and privacy  
 448 bound. As expected, smaller sparsity level  $q$  and smaller  $\sigma$  lead to larger privacy loss.

449 **4.2. Linear Regression.** We test the linear regression problem where the goal  
 450 is to minimize the objective  $f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \|\mathbf{y}_i - D_i \mathbf{x}\|^2$ , where  $D_i \in \mathbb{R}^{200 \times 5}$  and  
 451  $\mathbf{y}_i \in \mathbb{R}^{200}$  denote the local measurement matrix and local measurement vector at node  
 452  $i$ , respectively. To generate the data we synthesize the optimal solution  $\mathbf{x}^*$  from the  
 453 normal distribution. Then,  $\mathbf{y}_i$  is formed as  $\mathbf{y}_i = D_i \mathbf{x}^* + \xi_i$ , where  $\xi_i$  denotes the noise  
 454 added to the local measurement at node  $i$ . For all  $i$ ,  $D_i$  is drawn at random from  
 455 the standard normal distribution and then normalized so that its rows sum to 1;  $\xi_i$   
 456 is generated from a Gaussian distribution with zero mean and variance 0.01. Local  
 457 vectors  $\mathbf{x}_i^0$  are randomly initialized; the stepsize decreases with the iterations and is  
 458 set to  $\alpha_k = \frac{0.2}{k}$  in the  $k$ -th iteration.

459 In the implementation of the proposed algorithm, the gradient bound is set to  
 460  $D = 10$  and the privacy parameter  $\delta$  is set to  $\delta = 10^{-5}$ . We compute the residual for

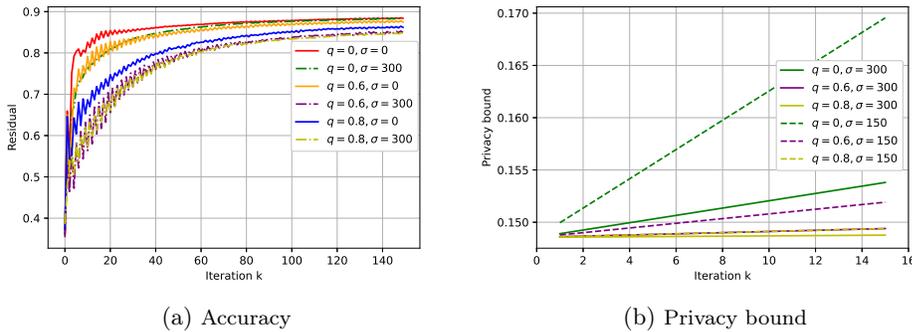


Fig. 3: Results on logistic regression on Stackoverflow. In (a) we observe sparsity and privacy delay convergence but they do not affect performance. In (b) we show the privacy loss over several iterations; we are able to maintain a reasonable budget for all combinations of parameters.

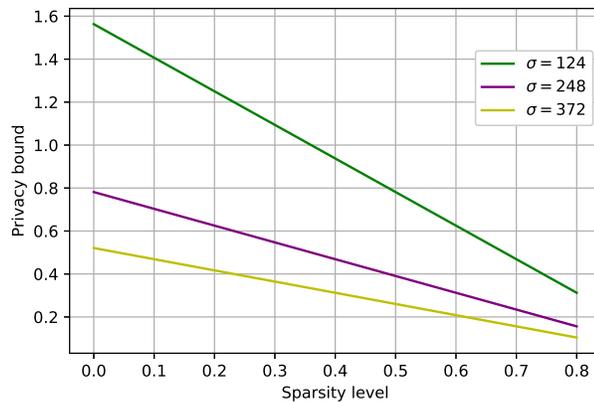


Fig. 4: Privacy bound for varied sparsity levels.

461 each iteration and show the results in Figure 5a. We observe that both sparsity and  
 462 privacy slow down the convergence.

463 Privacy bound for schemes with varied values of parameters are shown in Figure  
 464 5b, illustrating how privacy degrades over iterations.

### 465 4.3. Logistic Regression.

#### 466 4.3.1. Datasets.

467 *Stackoverflow*. Following prior work [41], we use a build vocabulary with 10000  
 468 frequent words and restrict each user’s dataset to have at most 128 sentences. We  
 469 rely on padding and truncation to enforce 20 word sentences, and represent them  
 470 with index sequences corresponding to the vocabulary words, out of vocabulary words,  
 471 beginning and end of sentences.

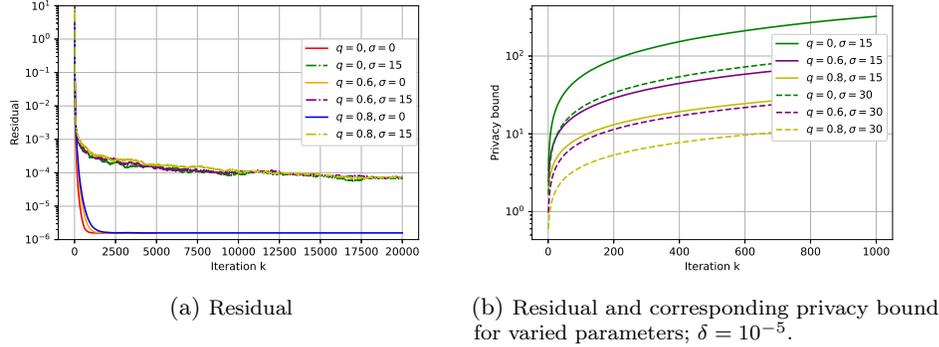


Fig. 5: Results of algorithm 2 on linear regression for synthetic data.  $B = 5$ . In (a) we show the sparsification and privacy will delay the convergence. In (b) we show how the privacy bound increases as we increase sparsity level and noise standard deviation.

472 The 150,000 data points are randomly split into 10 groups of equal size, where  
 473 each group is interpreted as being the local data for one of the nodes in the network.  
 474 Each node uses 13,500 data points as training data and the remaining 1500 points as  
 475 the validation set. The testing data contains 37640 data points.

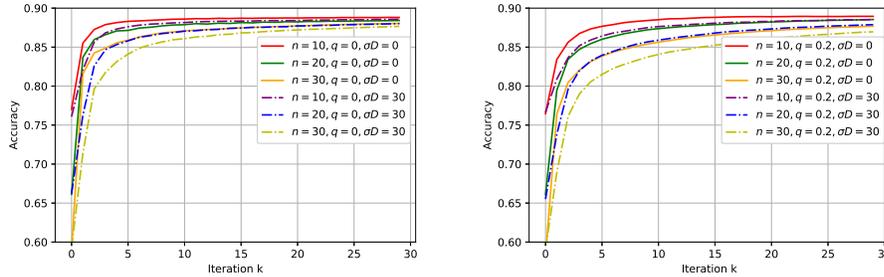


Fig. 6: Accuracy for varied network size.

476 **4.3.2. Varying the network size.** In this experiment we explore how the size of  
 477 the network affects the performance of the proposed algorithm 2. Here the total number  
 478 of data points is fixed and the number of local data points is inversely proportional  
 479 to the number of nodes in the network. In Figure 6, we see that the increasing the  
 480 network size and adding more noise delays the convergence without having much effect  
 481 on the final accuracy.

482 **4.3.3. Different topology.** So far, the network topology is randomly generated  
 483 at each iteration according to Erdős–Rényi model with some removing edges to make  
 484 the graph directed. Next, we consider a different type of the generative model. In  
 485 particular, we consider a topology periodically varying between the two networks  
 486 shown in Figure 7. This is a much sparser network than the previous ones, rendering  
 487 the algorithms slower as reflected by the results in Figure 5.

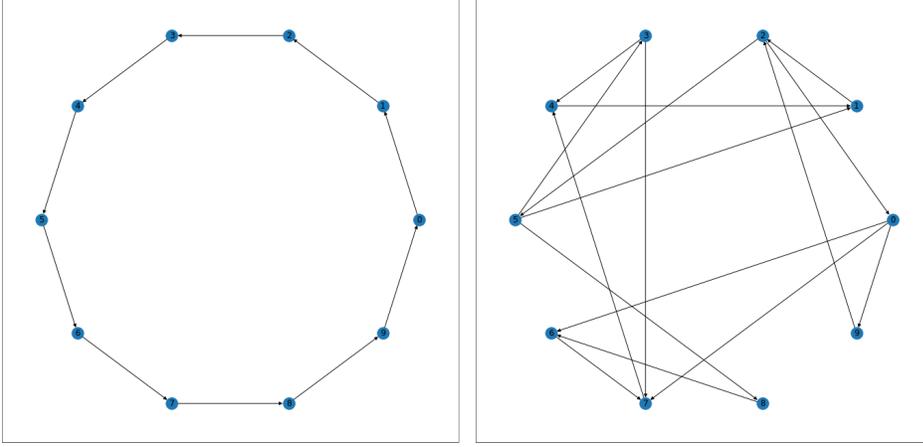


Fig. 7: Topology for periodically changing network

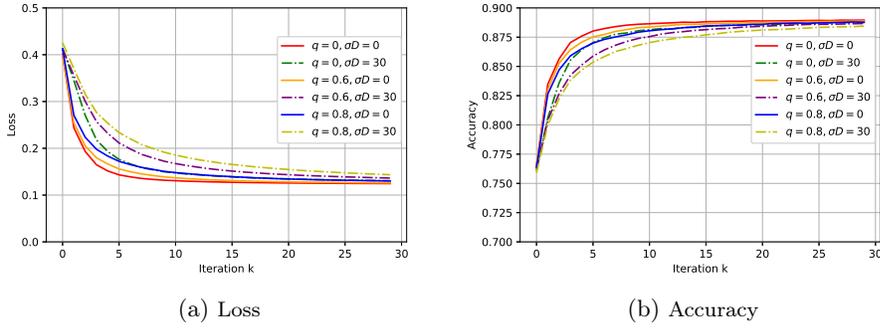


Fig. 8: Results of running the proposed algorithm 2 in the decentralized logistic regression model on a periodic network in Figure 4.

488 In particular, Figure 5 shows the loss and accuracy for the logistic regression  
 489 model of algorithm 2; all the remaining parameters of the experiment are the same as  
 490 in the main paper: the standard deviation  $\sigma D = 30$ , the step size  $\alpha_t = \frac{0.02}{t}$  and the  
 491 privacy parameter  $\delta = 10^{-4}$ .

492 Both privacy and sparsity constraints delay the convergence but do not affect the  
 493 final loss and accuracy.

#### 494 5. Formal convergence theorems and proofs.

THEOREM 2.4. (Theorem 2.4 in main body). Suppose Assumption 2.4 and 2.5 hold. Fix

$$\gamma \in \left(0, \min_m \left\{ \frac{1}{(20 + 8n)^n} (1 - |\lambda_3(\bar{M}_m((k+1)\mathcal{B} - 1 : k\mathcal{B}))|)^n \right\} \right),$$

495 and let  $\tau = \max_{C \in \mathcal{U}_M} |\lambda_2(C)| < 1$ ,  $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^0$ , and  $t \geq 0$ . Then running  
 496 Algorithm 2.1 for  $t$  iterations, suppose  $t = k\mathcal{B} - 1 + t'$ , where  $t' = 0, \dots, \mathcal{B} - 1$  and it

497 holds that for any  $i \in [n]$  and  $t \geq 1$ ,

$$498 \quad (5.1) \quad \begin{aligned} \|\mathbf{x}_i^t - \bar{\mathbf{z}}^t\| &\leq \sqrt{2nd}(\tau^{1/\mathcal{B}})^{t-(t'-1)} \sum_{j=1}^{2n} \sum_{m=1}^d |z_{jm}^0|, \\ \|\mathbf{y}_i^t\| &\leq \sqrt{2nd}(\tau^{1/\mathcal{B}})^{t-(t'-1)} \sum_{j=1}^{2n} \sum_{m=1}^d |z_{jm}^0|, \end{aligned}$$

499 where  $\bar{\mathbf{z}}^t = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^t + \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i^t$ . Further,  $E[\mathbf{x}_i^t]$  converges to  $\bar{\mathbf{x}}$  at a linear rate  
500  $O(\tau^{t/\mathcal{B}})$ .

501 *Proof.* To start with, we observe that the update in Algorithm 1 can be simplified  
502 as

$$503 \quad (5.2) \quad \begin{aligned} z_{im}^{t+1} &= \sum_{j=1}^{2n} [\bar{M}_m^t]_{ij} [Q(\mathbf{z}_j^t)]_m + \mathbb{1}_{\{t \bmod \mathcal{B} = \mathcal{B}-1\}} \gamma [F]_{ij} z_{jm}^{\mathcal{B}\lfloor t/\mathcal{B} \rfloor} \\ &= \sum_{j=1}^{2n} [\bar{M}_m^t]_{ij} z_{jm}^t + \mathbb{1}_{\{t \bmod \mathcal{B} = \mathcal{B}-1\}} \gamma [F]_{ij} z_{jm}^{\mathcal{B}\lfloor t/\mathcal{B} \rfloor} \end{aligned}$$

504 This holds because the mixing matrix is constructed such that its entries which multiply  
505 zero-valued (i.e., ‘‘sparsified’’) entries of  $Q(z_j^t)$  are set to be zero themselves. Next, we  
506 review the following lemma which help complete the proof after the incorporation of  
507 the noise expectation.

508 **LEMMA 5.1.** [*Theorem 2.4 in [10]*] Suppose Assumptions 2.4 and 2.5 hold, and  
509 instate the notations and hypotheses above. Then, there exist  $\sigma \in (0, 1)$  and  $\Gamma = \sqrt{2nd}$   
510 such that the following statements hold.

511 (a) For  $1 \leq i \leq n$  and  $t = k\mathcal{B} - 1 + t'$ , where  $t' = 0, \dots, \mathcal{B} - 1$ ,

$$512 \quad (5.3) \quad \|\mathbf{z}_i^t - \bar{\mathbf{z}}\| \leq \Gamma(\tau^{1/\mathcal{B}})^{t-(t'-1)} \sum_{j=1}^{2n} \sum_{m=1}^d |z_{jm}^0|,$$

513 where  $\bar{\mathbf{z}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^0 + \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i^0$ ;

514 (b) For  $1 + n \leq i \leq 2n$  and  $t = k\mathcal{B} - 1 + t'$ , where  $t' = 0, \dots, \mathcal{B} - 1$ ,

$$515 \quad (5.4) \quad \|\mathbf{z}_i^t\| \leq \Gamma(\tau^{1/\mathcal{B}})^{t-(t'-1)} \sum_{j=1}^{2n} \sum_{m=1}^d |z_{jm}^0|.$$

516 Now we can continue the proof of our theorem. In particular, using Lemma 5.1  
517 above, we have the first part (inequality) in the theorem proved, and establish  $\bar{\mathbf{z}}^t = \bar{\mathbf{z}}^0$   
518 for all  $t \geq 0$ . Since the noise added in the initialization part is unbiased, we have that  
519  $E[\bar{\mathbf{z}}^0] = \mathbf{z}^0$  where  $\mathbf{z}^0$  represents the initialization without noise. Then we can conclude  
520  $\lim_{t \rightarrow \infty} E[\mathbf{x}_i^t] = \bar{\mathbf{z}}^0 = \bar{\mathbf{x}}$  and the convergence rate is  $O(\tau^{t/\mathcal{B}})$ .  $\square$

521 **THEOREM 2.6.** (*Theorem 2.6 in main body*) Suppose Assumptions 2.4-2.5 on mix-  
522 ing matrices hold, fix  $\gamma \in (0, \min_m \left\{ \frac{1}{(20+8n)^n} (1 - |\lambda_3(\bar{M}_m((k+1)\mathcal{B} - 1 : k\mathcal{B}))|)^n \right\})$ .  
523 Assume  $|g_{im}^t| \leq D$ . If the objective function  $f$  is convex, and the step size  $\alpha_t$  decays  
524 as  $O(1/\sqrt{T})$ , Algorithm 2.2 converges at a rate of  $O(\ln T/\sqrt{T})$ , Concretely,

$$(5.5) \quad (E[f_{\min,T}] - f^*) \leq \frac{C_1}{\sum_{k=0}^{\lfloor T/\mathcal{B} \rfloor} \alpha_k} + \frac{C_2 \sum_{k=0}^{\lfloor T/\mathcal{B} \rfloor} \alpha_k^2}{\sum_{k=0}^{\lfloor T/\mathcal{B} \rfloor} \alpha_k} \leq \frac{C_1}{\sqrt{T/\mathcal{B} - 1}} + \frac{C_2 \ln(T/\mathcal{B})}{\sqrt{T/\mathcal{B} - 1}}$$

where  $f_{\min,T} := \min_{t=1, \dots, T} f(\mathbf{z}^t)$  and  $f^*$  is the optimal value, and

$$(5.6) \quad C_1 = \frac{n\sqrt{d}D^2}{2} + \frac{\sqrt{2}dD}{\sqrt{n}} \sum_{j=1}^{2n} \frac{\|\mathbf{z}_j^0\|}{1 - \tau^2},$$

528

$$(5.7) \quad C_2 = \frac{(d + \sigma^2)nD^2}{2} + \frac{\sqrt{2}dD}{\sqrt{n}} \sum_{j=1}^{2n} \|\mathbf{z}_j^0\| + \frac{2\sqrt{2nd^2(d + \sigma^2)}D^2}{1 - \tau} + \frac{4\sqrt{d(d + \sigma^2)}D^2}{n}$$

530 *Proof.* Similar to the steps in the consensus case, we start by the following  
531 observation to simplify the update in Algorithm 2:

$$\begin{aligned} z_{im}^{t+1} &= \sum_{j=1}^{2n} [\bar{M}_m^t]_{ij} [Q(\mathbf{z}_j^t)]_m + \mathbb{1}_{\{t \bmod \mathcal{B} = \mathcal{B} - 1\}} \gamma [F]_{ij} z_{jm}^{\mathcal{B} \lfloor t/\mathcal{B} \rfloor} \\ &\quad - \mathbb{1}_{\{t \bmod \mathcal{B} = \mathcal{B} - 1\}} \alpha_{\lfloor t/\mathcal{B} \rfloor} (g_{im}^{\mathcal{B} \lfloor t/\mathcal{B} \rfloor} + N(0, \sigma^2)) \\ &= \sum_{j=1}^{2n} [\bar{M}_m^t]_{ij} z_{jm}^t + \mathbb{1}_{\{t \bmod \mathcal{B} = \mathcal{B} - 1\}} \gamma [F]_{ij} z_{jm}^{\mathcal{B} \lfloor t/\mathcal{B} \rfloor} \\ &\quad - \mathbb{1}_{\{t \bmod \mathcal{B} = \mathcal{B} - 1\}} \alpha_{\lfloor t/\mathcal{B} \rfloor} (g_{im}^{\mathcal{B} \lfloor t/\mathcal{B} \rfloor} + N(0, \sigma^2)) \end{aligned}$$

533 Let  $[\nabla \tilde{f}_i(\mathbf{z}_i^t)]_m = g_{im}^t + N(0, \sigma^2 D^2) \mathbb{1}_{\{i \leq n\}}$  be the  $m^{\text{th}}$  entry of  $\nabla \tilde{f}_i(\mathbf{z}_i^t)$  and we  
534 then compute

$$(5.8) \quad \|\bar{\mathbf{z}}^{(k+1)\mathcal{B}} - \mathbf{x}^*\|^2 = \|\bar{\mathbf{z}}^{k\mathcal{B}} - \mathbf{x}^*\|^2 + \left\| \frac{\alpha_k}{n} \sum_{i=1}^n \nabla \tilde{f}_i(\mathbf{z}_i^{k\mathcal{B}}) \right\|^2 - 2 \frac{\alpha_k}{n} \sum_{i=1}^n \langle \bar{\mathbf{z}}^{k\mathcal{B}} - \mathbf{x}^*, \nabla \tilde{f}_i(\mathbf{z}_i^{k\mathcal{B}}) \rangle.$$

536

$$\begin{aligned} E[\|\bar{\mathbf{z}}^{(k+1)\mathcal{B}} - \mathbf{x}^*\|^2 | \mathcal{F}_{k\mathcal{B}}] &= \|\bar{\mathbf{z}}^{k\mathcal{B}} - \mathbf{x}^*\|^2 + E\left[\left\| \frac{\alpha_k}{n} \sum_{i=1}^n \nabla \tilde{f}_i(\mathbf{z}_i^{k\mathcal{B}}) \right\|^2 | \mathcal{F}_{k\mathcal{B}}\right] \\ &\quad - 2 \frac{\alpha_k}{n} \sum_{i=1}^n \langle \bar{\mathbf{z}}^{k\mathcal{B}} - \mathbf{x}^*, \nabla f_i(\mathbf{z}_i^{k\mathcal{B}}) \rangle \\ &\leq \|\bar{\mathbf{z}}^{k\mathcal{B}} - \mathbf{x}^*\|^2 + E\left[\left\| \frac{\alpha_k}{n} \sum_{i=1}^n \nabla \tilde{f}_i(\mathbf{z}_i^{k\mathcal{B}}) \right\|^2 | \mathcal{F}_{k\mathcal{B}}\right] \\ &\quad - 2 \frac{\alpha_k}{n} \sum_{i=1}^n (-2\sqrt{d}D \|\bar{\mathbf{z}}^{k\mathcal{B}} - \mathbf{z}_i^{k\mathcal{B}}\| + f_i(\bar{\mathbf{z}}^{k\mathcal{B}}) - f_i(\mathbf{x}^*)) \end{aligned}$$

542 where the last inequality is derived from

$$\begin{aligned} \langle \bar{\mathbf{z}}^{k\mathcal{B}} - \mathbf{x}^*, \nabla f_i(\mathbf{z}_i^{k\mathcal{B}}) \rangle &\geq \langle \bar{\mathbf{z}}^{k\mathcal{B}} - \mathbf{z}_i^{k\mathcal{B}}, \nabla f_i(\mathbf{z}_i^{k\mathcal{B}}) \rangle + f_i(\mathbf{z}_i^{k\mathcal{B}}) - f_i(\mathbf{x}^*) \\ &\geq -\sqrt{d}D \|\bar{\mathbf{z}}^{k\mathcal{B}} - \mathbf{z}_i^{k\mathcal{B}}\| + f_i(\mathbf{z}_i^{k\mathcal{B}}) - f_i(\bar{\mathbf{z}}^{k\mathcal{B}}) + f_i(\bar{\mathbf{z}}^{k\mathcal{B}}) - f_i(\mathbf{x}^*) \\ &\geq -2\sqrt{d}D \|\bar{\mathbf{z}}^{k\mathcal{B}} - \mathbf{z}_i^{k\mathcal{B}}\| + f_i(\bar{\mathbf{z}}^{k\mathcal{B}}) - f_i(\mathbf{x}^*) \end{aligned}$$

546

547 Now, the unconditional expectation satisfies

$$548 \quad E[\|\bar{\mathbf{z}}^{(k+1)\mathcal{B}} - \mathbf{x}^*\|^2] \leq E[\|\bar{\mathbf{z}}^{k\mathcal{B}} - \mathbf{x}^*\|^2] + E\left[\left\|\frac{\alpha_k}{n} \sum_{i=1}^n \nabla \tilde{f}_i(\mathbf{z}_i^{k\mathcal{B}})\right\|^2\right] + 4\frac{\alpha_k \sqrt{d}D}{n} \sum_{i=1}^n E[\|\bar{\mathbf{z}}^{k\mathcal{B}} - \mathbf{z}_i^{k\mathcal{B}}\|]$$

$$549 \quad - 2\frac{\alpha_k}{n} \sum_{i=1}^n E[f_i(\bar{\mathbf{z}}^{k\mathcal{B}})] - f_i(x^*).$$

550

551 Summing over  $k$  from 0 to  $\infty$  and rearranging yields

$$552 \quad (5.9) \quad 2 \sum_{k=0}^{\lfloor T/\mathcal{B} \rfloor} \alpha_k (E[f(\bar{\mathbf{z}}^{k\mathcal{B}})] - f^*) \leq \|\bar{\mathbf{z}}^0 - \mathbf{x}^*\|^2 + n(d + \sigma^2)D^2 \sum_{k=0}^{\lfloor T/\mathcal{B} \rfloor} \alpha_k^2$$

$$+ \frac{4\sqrt{d}D}{n} \sum_{i=1}^n \sum_{k=0}^{\lfloor T/\mathcal{B} \rfloor} \alpha_k E[\|\bar{\mathbf{z}}^{k\mathcal{B}} - \mathbf{z}_i^{k\mathcal{B}}\|].$$

553 Similar to the derivations for Lemma 3 in [10], we next obtain an upper bound for  
554 the last term in (5.9).

555 Since the update in Algorithm 2.2 implies

$$556 \quad z_{im}^{k\mathcal{B}} = \sum_{j=1}^n [M_m(k\mathcal{B} - 1 : 0)]_{ij} z_{jm}^0 - \sum_{r=1}^{k-1} \sum_{j=1}^{2n} [M_m((k-1)\mathcal{B} - 1 : (r-1)\mathcal{B})]_{ij} \alpha_{r-1} [\nabla \tilde{f}_j(\mathbf{z}_j^{(r-1)\mathcal{B}})]_m$$

$$- \alpha_{k-1} [\nabla \tilde{f}_i(\mathbf{z}_i^{(k-1)\mathcal{B}})]_m$$

557 for  $i \in \{1, \dots, 2n\}$  and  $m \in \{1, \dots, d\}$  and using the fact that the mixing matrix and  
558 its product have column sum equal to 1,

$$559 \quad \bar{z}_m^{k\mathcal{B}} = \frac{1}{n} \sum_{j=1}^{2n} z_{jm}^{k\mathcal{B}}$$

$$= \frac{1}{n} \sum_{j=1}^{2n} z_{jm}^0 - \frac{1}{n} \sum_{r=1}^{(k-1)\mathcal{B}} \sum_{j=1}^{2n} [\nabla \tilde{f}_j(\mathbf{z}_j^{(r-1)\mathcal{B}})]_m - \frac{1}{n} \sum_{j=1}^{2n} \alpha_{k-1} [\nabla \tilde{f}_j(\mathbf{z}_j^{(k-1)\mathcal{B}})]_m$$

560 Then using the gradient norm bound and noise variance, we derive the following upper  
561 bound for the last term in (5.9)

$$562 \quad (5.10) \quad \sum_{i=1}^n \sum_{k=0}^{\lfloor T/\mathcal{B} \rfloor} \alpha_k E[\|\bar{\mathbf{z}}^{k\mathcal{B}} - \mathbf{z}_i^{k\mathcal{B}}\|] \leq \sqrt{2nd} \sum_{j=1}^{2n} \|\mathbf{z}_j^0\| \sum_{k=1}^{\lfloor T/\mathcal{B} \rfloor} \alpha_k \tau^k$$

$$+ \sqrt{2nd(d + \sigma^2)} nD \sum_{k=1}^{\lfloor T/\mathcal{B} \rfloor} \sum_{r=1}^{k-1} \tau^{k-r} \alpha_k \alpha_{r-1}$$

$$+ 2\sqrt{d + \sigma^2} D \sum_{k=0}^{\lfloor T/\mathcal{B} \rfloor - 1} \alpha_k^2$$

563 Applying  $ab \leq \frac{1}{2}(a + b)^2$ , we have following bounds:

$$564 \quad (5.11) \quad \sum_{k=1}^{\lfloor T/\mathcal{B} \rfloor} \alpha_k \tau^k \leq \frac{1}{2} \sum_{k=1}^{\lfloor T/\mathcal{B} \rfloor} \alpha_k^2 + \frac{1}{1 - \tau^2}$$

565 (5.12)

$$566 \sum_{k=1}^{\lfloor T/\mathcal{B} \rfloor} \sum_{r=1}^{k-1} \tau^{k-r} \alpha_k \alpha_{r-1} \leq \frac{1}{2} \sum_{k=1}^{\lfloor T/\mathcal{B} \rfloor} \alpha_k^2 \sum_{r=1}^{k-1} \tau^{k-r} + \frac{1}{2} \sum_{r=1}^{\lfloor T/\mathcal{B} \rfloor - 1} \alpha_{r-1}^2 \sum_{k=r+1}^{\lfloor T/\mathcal{B} \rfloor - 1} \tau^{k-r} \leq \frac{1}{1-\tau} \sum_{k=1}^{\lfloor T/\mathcal{B} \rfloor} \alpha_k^2$$

567 Then

$$568 E[\min_{t=1, \dots, T} f(\bar{\mathbf{z}}^t)] \rightarrow \frac{1}{n} \sum_{i=1}^n f_i(x^*) = f^*.$$

569 Defining  $f_{\min} := \min_t f(\bar{\mathbf{z}}^t)$ , we have

$$570 (5.13) \quad (E[f_{\min}] - f^*) \sum_{k=0}^{\lfloor T/\mathcal{B} \rfloor} \alpha_k \leq \sum_{k=0}^{\lfloor T/\mathcal{B} \rfloor} \alpha_k (f(\bar{\mathbf{z}}^{k\mathcal{B}}) - f^*) \leq C_1 + C_2 \sum_{k=0}^{\lfloor T/\mathcal{B} \rfloor} \alpha_k^2,$$

571 where

$$572 (5.14) \quad C_1 = \frac{n\sqrt{d}D^2}{2} + \frac{\sqrt{2}dD}{\sqrt{n}} \sum_{j=1}^{2n} \frac{\|\mathbf{z}_j^0\|}{1-\tau^2},$$

573

$$574 (5.15) \quad C_2 = \frac{(d+\sigma^2)nD^2}{2} + \frac{\sqrt{2}dD}{\sqrt{n}} \sum_{j=1}^{2n} \|\mathbf{z}_j^0\| + \frac{2\sqrt{2nd^2(d+\sigma^2)}D^2}{1-\tau} + \frac{4\sqrt{d(d+\sigma^2)}D^2}{n}$$

575 Note that we can express (5.13) equivalently as

$$576 (5.16) \quad (E[f_{\min}] - f^*) \leq \frac{C_1}{\sum_{k=0}^{\lfloor T/\mathcal{B} \rfloor} \alpha_k} + \frac{C_2 \sum_{k=0}^{\lfloor T/\mathcal{B} \rfloor} \alpha_k^2}{\sum_{k=0}^{\lfloor T/\mathcal{B} \rfloor} \alpha_k}.$$

577 If we select the schedule of stepsizes according to  $\alpha_t = O(1/\sqrt{t})$ , the first term on the  
578 right hand side of (5.16) satisfies

$$579 (5.17) \quad \frac{C_1}{\sum_{t=0}^T \alpha_t} = C_1 \frac{1/2}{\sqrt{T} - 1}$$

580 while for the second term it holds that

$$581 (5.18) \quad \frac{C_2 \sum_{t=0}^T \alpha_t^2}{\sum_{t=0}^T \alpha_t} = C_2 \frac{\ln T}{2(\sqrt{T} - 1)}$$

582 Recall that  $\sigma = \mathcal{O}(\frac{\sqrt{T} \ln(1/\delta)}{\epsilon})$ , then

$$583 C_2 \frac{\ln T}{2(\sqrt{T} - 1)} = \mathcal{O}\left(\frac{\sqrt{d}}{1-\tau} \frac{\ln T}{(\sqrt{T} - 1)} + \left(d + \frac{T(\ln(1/\delta))^2}{\epsilon^2}\right) \frac{\ln T}{(\sqrt{T} - 1)}\right) \quad \square$$

584 **THEOREM 2.8.** (Theorem 2.8 in main body) Suppose Assumptions 2.4-2.5 on mix-  
585 ing matrices hold, fix  $\gamma \in (0, \min_m \left\{ \frac{1}{(20+8n)^n} (1 - |\lambda_3(\bar{M}_m((k+1)\mathcal{B} - 1 : k\mathcal{B})|)|)^n \right\})$ .586 Assume  $|g_{im}^t| \leq D$  and let  $\hat{\mathbf{x}}_i^T = \frac{\sum_{t=1}^{\lfloor T/\mathcal{B} \rfloor} (t-1)\mathbf{x}_i^t}{t(t-1)/2}$  for  $T \geq \mathcal{B}$  and  $K = \lfloor T/\mathcal{B} \rfloor$ . If the

587 objective function  $f$  is convex and the local objective function  $f_i$  is  $\mu_i$  strongly-convex,  
588 then there exists some constant  $C_3 > 0, C_4 > 0$  such that for all  $i$

(5.19)

$$E[f(\hat{\mathbf{x}}_i^T) - f(\mathbf{x}^*)] \leq \frac{C_3}{K} \sum_{j=1}^n \|\mathbf{x}_j^0\|_1 + \frac{C_4}{K} (1 + \ln(K-1)) + \frac{p^2 n D^2 d}{K} (1 + \sigma^2)$$

589

$$E\left[\sum_{j=1}^n \mu_j \|\hat{\mathbf{x}}_j^T - \mathbf{x}^*\|^2\right] \leq \frac{C_3}{K} \sum_{j=1}^n \|\mathbf{x}_j^0\|_1 + \frac{C_4}{K} (1 + \ln(K-1)) + \frac{p^2 n D^2 d}{K} (1 + \sigma^2)$$

590 where the step size  $\alpha_t = \frac{p}{t}$  for  $p \geq \frac{4n}{\sum_{i=1}^n \mu_i}$ ,  $\mathbf{x}^*$  is the optimal solution,  $C_3 = \frac{5\sqrt{2nndD}}{1-\tau}$

591 and  $C_4 = \frac{5\sqrt{2nd(1+\sigma^2)n^2dD^2}}{1-\tau}$ .

592 *Proof.* Let  $\mathbf{v} \in \mathcal{R}^d$  be any arbitrary vector and  $\xi_j^t$  be the noise vector at  $\nabla f_j(\mathbf{x}_j^t)$   
593 (denoted as  $\nabla f_j^t$  shortly). Since  $\mathbf{g}_j^t = \nabla f_j(\mathbf{x}_j^t) = \nabla f_j^t$  when  $i \in \{1, \dots, n\}$ , for all  
594  $t \geq 0$ ,

$$\begin{aligned} z_{im}^{t+1} &= \sum_{j=1}^{2n} [\bar{M}_m^t]_{ij} [Q(\mathbf{z}_j^t)]_m + \mathbb{1}_{\{t \bmod \mathcal{B}=\mathcal{B}-1\}} \gamma [F]_{ij} z_{jm}^{\mathcal{B}[t/\mathcal{B}]} \\ &\quad - \mathbb{1}_{\{t \bmod \mathcal{B}=\mathcal{B}-1\}} \alpha_{[t/\mathcal{B}]} (g_{im}^{\mathcal{B}[t/\mathcal{B}]} + \xi_{jm}^{\mathcal{B}[t/\mathcal{B}]} \mathbb{1}_{(i \leq n)}) \\ &= \sum_{j=1}^{2n} [\bar{M}_m^t]_{ij} z_{jm}^t + \mathbb{1}_{\{t \bmod \mathcal{B}=\mathcal{B}-1\}} \gamma [F]_{ij} z_{jm}^{\mathcal{B}[t/\mathcal{B}]} \\ &\quad - \mathbb{1}_{\{t \bmod \mathcal{B}=\mathcal{B}-1\}} \alpha_{[t/\mathcal{B}]} (g_{im}^{\mathcal{B}[t/\mathcal{B}]} + \xi_{jm}^{\mathcal{B}[t/\mathcal{B}]} \mathbb{1}_{(i \leq n)}) \end{aligned}$$

595

596 and, moreover,

$$\begin{aligned} \bar{z}_m^{t+1} &= \frac{1}{n} \sum_{j=1}^{2n} z_{jm}^t - \frac{1}{n} \alpha_{[t/\mathcal{B}]} \sum_{j=1}^{2n} \mathbb{1}_{\{t \bmod \mathcal{B}=\mathcal{B}-1\}} (g_{im}^{\mathcal{B}[t/\mathcal{B}]} + \xi_{jm}^{\mathcal{B}[t/\mathcal{B}]} \mathbb{1}_{(i \leq n)}) \\ &= \bar{z}_m^t - \frac{\alpha_{[t/\mathcal{B}]}}{n} \sum_{j=1}^{2n} (g_{jm}^{\mathcal{B}[t/\mathcal{B}]} + \xi_j^{\mathcal{B}[t/\mathcal{B}]} \mathbb{1}_{(j \leq n)}). \end{aligned}$$

597 (5.20)

598 Then,

(5.21)

$$\|\bar{\mathbf{z}}^{(k+1)\mathcal{B}} - \mathbf{v}\|^2 = \|\bar{\mathbf{z}}^{k\mathcal{B}} - \mathbf{v}\|^2 - \frac{\alpha_k}{n} \sum_{j=1}^{2n} (\mathbf{g}_j^{k\mathcal{B}} + \xi_j^{k\mathcal{B}} \mathbb{1}_{(j \leq n)})' (\bar{\mathbf{z}}^{k\mathcal{B}} - \mathbf{v}) + \frac{\alpha_k^2}{n^2} \left\| \sum_{j=1}^{2n} \mathbf{g}_j^{k\mathcal{B}} + \xi_j^{k\mathcal{B}} \mathbb{1}_{(j \leq n)} \right\|^2$$

(5.22)

$$= \|\bar{\mathbf{z}}^{k\mathcal{B}} - \mathbf{v}\|^2 - \frac{\alpha_k}{n} \sum_{j=1}^n (\nabla f_j^{k\mathcal{B}} + \xi_j^{k\mathcal{B}})' (\bar{\mathbf{z}}^{k\mathcal{B}} - \mathbf{v}) + \frac{\alpha_k^2}{n^2} \left\| \sum_{j=1}^n \nabla f_j^{k\mathcal{B}} + \xi_j^{k\mathcal{B}} \right\|^2. \quad \blacksquare$$

600

601

602 Now, we rewrite each cross-term  $(\nabla f_j^t)'(\bar{\mathbf{z}}^t - \mathbf{v})$  as

$$\begin{aligned}
603 \quad (5.23) \quad (\nabla f_j^t)'(\bar{\mathbf{z}}^t - \mathbf{v}) &= (\nabla f_j^t)'(\bar{\mathbf{z}}^t - \mathbf{x}_j^t) + (\nabla f_j^t)'(\mathbf{x}_j^t - \mathbf{v}) \\
&\geq -\sqrt{d}D\|\bar{\mathbf{z}}^t - \mathbf{x}_j^t\| + f_j(\mathbf{x}_j^t) - f_j(\mathbf{v}) + \frac{\mu_j}{2}\|\mathbf{x}_j^t - \mathbf{v}\|^2 \\
&= -\sqrt{d}D\|\bar{\mathbf{z}}^t - \mathbf{x}_j^t\| + (f_j(\mathbf{x}_j^t) - f_j(\bar{\mathbf{z}}^t)) + (f_j(\bar{\mathbf{z}}^t) - f_j(\mathbf{v})) \\
&\quad + \frac{\mu_j}{2}\|\mathbf{x}_j^t - \mathbf{v}\|^2 \\
&\geq -2\sqrt{d}D\|\bar{\mathbf{z}}^t - \mathbf{x}_j^t\| + (f_j(\bar{\mathbf{z}}^t) - f_j(\mathbf{v})) + \frac{\mu_j}{2}\|\mathbf{x}_j^t - \mathbf{v}\|^2.
\end{aligned}$$

604 Using  $f(\mathbf{x}) = \frac{1}{n} \sum_{j=1}^n f_j(\mathbf{x})$ ,

$$605 \quad (5.24) \quad \sum_{j=1}^n (\nabla f_j^{k\mathcal{B}})'(\bar{\mathbf{z}}^{k\mathcal{B}} - \mathbf{v}) \geq n(f(\mathbf{x}^{k\mathcal{B}}) - f(\mathbf{v})) + \frac{1}{2} \sum_{j=1}^n \mu_j \|\mathbf{x}_j^{k\mathcal{B}} - \mathbf{v}\|^2 - 2 \sum_{j=1}^n \sqrt{d}D \|\bar{\mathbf{z}}^{k\mathcal{B}} - \mathbf{x}_j^{k\mathcal{B}}\|.$$

606 Hence, we have shown that

$$\begin{aligned}
607 \quad E[\|\bar{\mathbf{z}}^{(k+1)\mathcal{B}} - \mathbf{v}\|^2 | \mathcal{F}_{k\mathcal{B}}] &\leq \|\bar{\mathbf{z}}^{k\mathcal{B}} - \mathbf{v}\|^2 - 2\alpha_k(f(\bar{\mathbf{z}}^{k\mathcal{B}}) - f(\mathbf{v})) - \frac{\alpha_k}{n} \sum_{j=1}^n \mu_j \|\mathbf{x}_j^{k\mathcal{B}} - \mathbf{v}\|^2 \\
608 \quad &\quad + \frac{4\alpha_k}{n} \sum_{j=1}^n \sqrt{d}D \|\mathbf{x}_j^{k\mathcal{B}} - \bar{\mathbf{z}}^{k\mathcal{B}}\| + \frac{\alpha_k^2}{n^2} \sum_{j=1}^n (\sqrt{d}D + \sigma\sqrt{d}D)^2 \\
609 \quad &\leq \|\bar{\mathbf{z}}^{k\mathcal{B}} - \mathbf{v}\|^2 - 2\alpha_k(f(\bar{\mathbf{z}}^{k\mathcal{B}}) - f(\mathbf{v})) - \frac{\alpha_k}{n} \sum_{j=1}^n \mu_j \|\mathbf{x}_j^{k\mathcal{B}} - \mathbf{v}\|^2 \\
610 \quad &\quad + \frac{4\alpha_k}{n} \sum_{j=1}^n \sqrt{d}D \|\mathbf{x}_j^{k\mathcal{B}} - \bar{\mathbf{z}}^{k\mathcal{B}}\| + \frac{\alpha_k^2}{n^2} \sum_{j=1}^n d(1 + \sigma^2)D^2. \\
611
\end{aligned}$$

612 Then we can replace  $\mathbf{v}$  by the optimal solution  $\mathbf{x}^*$  and this gives

$$\begin{aligned}
613 \quad (5.25) \quad E[\|\bar{\mathbf{z}}^{(k+1)\mathcal{B}} - \mathbf{x}^*\|^2 | \mathcal{F}_{k\mathcal{B}}] &\leq \|\bar{\mathbf{z}}^{k\mathcal{B}} - \mathbf{x}^*\|^2 - 2\alpha_k(f(\bar{\mathbf{z}}^{k\mathcal{B}}) - f(\mathbf{x}^*)) - \frac{\alpha_k}{n} \sum_{j=1}^n \mu_j \|\mathbf{x}_j^{k\mathcal{B}} - \mathbf{x}^*\|^2 \\
614 \quad &\quad + \frac{4\alpha_k}{n} \sum_{j=1}^n \sqrt{d}D \|\mathbf{x}_j^{k\mathcal{B}} - \bar{\mathbf{z}}^{k\mathcal{B}}\| + \frac{\alpha_k^2}{n^2} \sum_{j=1}^n d(1 + \sigma^2)D^2. \quad \blacksquare
\end{aligned}$$

615 Since  $f(\mathbf{x}) = \frac{1}{n} \sum_{j=1}^n f_j(\mathbf{x})$  is convex, each local objective function  $f_i$  is  $\mu_i$ -strongly-  
616 convex and the upper bound on the gradient magnitude is  $|g_{im}^t| \leq D$ , the following  
617 two inequalities hold:

(a)

$$618 \quad (5.26) \quad f(\bar{\mathbf{z}}^t) - f(\mathbf{x}^*) \geq \frac{1}{2n} \left( \sum_{j=1}^n \mu_j \right) \|\bar{\mathbf{z}}^t - \mathbf{x}^*\|$$

(b)

$$619 \quad (5.27) \quad f(\bar{\mathbf{z}}^t) - f(\mathbf{x}^*) \geq -\frac{L}{n} \|\mathbf{x}_i^t - \bar{\mathbf{z}}^t\| + f(\mathbf{x}_i^t) - f(\mathbf{x}^*)$$

620 where  $L = n\sqrt{d}D$ , for any  $i = 1, \dots, n$ .

621 The above imply that for all  $i = 1, \dots, n$ ,

$$622 \quad (5.28) \quad 2(f(\bar{\mathbf{z}}^t) - f(\mathbf{x}^*)) \geq \frac{1}{2} \left( \frac{1}{n} \sum_{j=1}^n \mu_j \right) \|\bar{\mathbf{z}}^t - \mathbf{x}^*\| - \frac{L}{n} \|\mathbf{x}_i^t - \bar{\mathbf{z}}^t\| + f(\mathbf{x}_i^t) - f(\mathbf{x}^*).$$

623 Now, for each  $i = 1, \dots, n$ ,

$$624 \quad E[\|\bar{\mathbf{z}}^{(k+1)\mathcal{B}} - \mathbf{x}^*\|^2 | \mathcal{F}_{k\mathcal{B}}] \leq \|\bar{\mathbf{z}}^{k\mathcal{B}} - \mathbf{x}^*\|^2 - \frac{\alpha_k}{n} \left( \frac{1}{2} \left( \sum_{j=1}^n \mu_j \right) \|\bar{\mathbf{z}}^{k\mathcal{B}} - \mathbf{x}^*\| - L \|\mathbf{x}_i^{k\mathcal{B}} - \bar{\mathbf{z}}^{k\mathcal{B}}\| + n(f(\mathbf{x}_i^{k\mathcal{B}}) - f(\mathbf{x}^*)) \right) \\ 625 \quad - \frac{\alpha_k}{n} \sum_{j=1}^n \mu_j \|\mathbf{x}_j^{k\mathcal{B}} - \mathbf{x}^*\|^2 + \frac{4\alpha_k}{n} \sum_{j=1}^n \sqrt{d}D \|\mathbf{x}_j^{k\mathcal{B}} - \bar{\mathbf{z}}^{k\mathcal{B}}\| + \frac{\alpha_k^2}{n} \sum_{j=1}^n (\sqrt{d}D + \sigma\sqrt{d}D)^2. \quad \blacksquare$$

627 Let  $\alpha_k = \frac{p}{k+1}$ ; since  $p \frac{\sum_{i=1}^n \mu_i}{n} \geq 4$ ,

$$628 \quad E[\|\bar{\mathbf{z}}^{(k+1)\mathcal{B}} - \mathbf{x}^*\|^2 | \mathcal{F}_{k\mathcal{B}}] \leq \left(1 - \frac{2}{k+1}\right) \|\bar{\mathbf{z}}^{k\mathcal{B}} - \mathbf{x}^*\|^2 - \frac{p}{(k+1)} (f(\mathbf{x}_i^{k\mathcal{B}}) - f(\mathbf{x}^*)) \\ 629 \quad + \frac{pL}{n(k+1)} \|\mathbf{x}_i^{k\mathcal{B}} - \bar{\mathbf{z}}^{k\mathcal{B}}\| - \frac{p}{n(k+1)} \sum_{j=1}^n \mu_j \|\mathbf{x}_j^{k\mathcal{B}} - \mathbf{x}^*\|^2 \\ 630 \quad + \frac{4p}{n(k+1)} \sum_{j=1}^n \sqrt{d}D \|\mathbf{x}_j^{k\mathcal{B}} - \bar{\mathbf{z}}^{k\mathcal{B}}\| + \frac{p^2}{n(k+1)^2} \sum_{j=1}^n (\sqrt{d}D + \sigma\sqrt{d}D)^2. \quad \blacksquare$$

632 Multiply both sides of the above inequality by  $k(k+1)$  and taking the expectation  
633 yields for all  $T \geq \mathcal{B}$ , let  $K = \lfloor T/\mathcal{B} \rfloor$ ,

(5.29)

$$634 \quad K(K-1)E[\|\bar{\mathbf{z}}^{K\mathcal{B}} - \mathbf{x}^*\|^2] \leq -\frac{p}{n} \sum_{k=1}^{K-1} t E[n(f(\mathbf{x}_i^{k\mathcal{B}}) - f(\mathbf{x}^*)) + \sum_{j=1}^n \mu_j \|\mathbf{x}_j^{k\mathcal{B}} - \mathbf{x}^*\|^2] \\ 635 \quad + \frac{pL}{n} \sum_{k=1}^{K-1} k E[\|\mathbf{x}_i^{k\mathcal{B}} - \bar{\mathbf{z}}^{k\mathcal{B}}\|] - \frac{p}{n(k+1)} \left( \sum_{j=1}^n \mu_j \right) E[\|\mathbf{x}_j^{k\mathcal{B}} - \mathbf{x}^*\|^2] \\ 636 \quad + \frac{4p}{n} \sum_{k=1}^{K-1} t \sum_{j=1}^n \sqrt{d}D E[\|\mathbf{x}_j^{k\mathcal{B}} - \bar{\mathbf{z}}^{k\mathcal{B}}\|] + \frac{p^2}{n} \sum_{j=1}^n (\sqrt{d}D + \sigma\sqrt{d}D)^2 \sum_{k=1}^{K-1} \frac{k}{k+1}. \quad \blacksquare$$

636 To derive the upper bound on  $E[\sum_{k=1}^{K-1} \|\mathbf{x}_i^{k\mathcal{B}} - \bar{\mathbf{z}}^{k\mathcal{B}}\|]$ , we refer to Lemma 3 in [10]  
637 and Corollary 1 and 2 in [36]. In particular, we have the following

$$638 \quad (5.30) \quad E\left[\sum_{k=1}^K \|\mathbf{x}_i^{k\mathcal{B}} - \bar{\mathbf{z}}^{k\mathcal{B}}\|\right] \leq C'_3 \sum_{j=1}^n \|\mathbf{x}_j^0\|_1 + C'_4(1 + \ln K).$$

639 where

$$640 \quad (5.31) \quad C'_3 = \frac{\sqrt{2nd}}{1-\tau}, \quad C'_4 = \frac{\sqrt{2nd^2(1+\sigma^2)}nD}{1-\tau}$$

641 Therefore,

(5.32)

$$\begin{aligned} \frac{1}{K(K-1)} \sum_{k=1}^{K-1} tE[n(f(\mathbf{x}_i^{kB}) - f(\mathbf{x}^*)) + \sum_{j=1}^n \mu_j \|\mathbf{x}_j^{kB} - \mathbf{x}^*\|^2] &\leq \frac{5L}{T} (C'_3 \sum_{j=1}^n \|\mathbf{x}_j^0\|_1 \\ &+ C'_4(1 + \ln(K-1))) \\ &+ \frac{p^2}{K} \sum_{j=1}^n d(1 + \sigma^2)D^2 \end{aligned}$$

642

643

644 Hence,

(5.33)

$$\begin{aligned} \frac{1}{K(K-1)} \sum_{k=1}^{K-1} kE[(f(\mathbf{x}_i^{kB}) - f(\mathbf{x}^*)) + \sum_{j=1}^n \mu_j \|\mathbf{x}_j^{kB} - \mathbf{x}^*\|^2] &\leq \frac{5L}{K} (C'_3 \sum_{j=1}^n \|\mathbf{x}_j^0\|_1 \\ &+ C'_4(1 + \ln(K-1))) + \\ &\frac{p^2}{K} \sum_{j=1}^n d(1 + \sigma^2)D^2 \end{aligned}$$

645

646

By convexity, for each  $i \in [n]$  it holds that

(5.34)

$$\frac{2}{K(K-1)} \sum_{k=1}^{K-1} t(f(\mathbf{x}_i^{kB}) - f(\mathbf{x}^*)) + \sum_{j=1}^n \mu_j \|\mathbf{x}_j^{kB} - \mathbf{x}^*\|^2 \geq f(\hat{\mathbf{x}}^K) - f(\mathbf{x}^*) + \sum_{j=1}^n \mu_j \|\hat{\mathbf{x}}^K - \mathbf{x}^*\|^2,$$

648

649 where  $\hat{\mathbf{x}}_i^K = \frac{\sum_{k=1}^K (k-1)\mathbf{x}_i^k}{k(k-1)/2}$  for  $K \geq 2$ .

650 Setting  $C_3 = 5LC'_3 = \frac{5\sqrt{2nnd}D}{1-\tau}$  and  $C_4 = 5LC'_4 = \frac{5\sqrt{2nd(1+\sigma^2)n^2dD^2}}{1-\tau}$  completes  
651 the proof.

652 **6. Conclusion.** In this paper we propose differentially private and communica-  
653 tion efficient algorithms for decentralized consensus and optimization over directed  
654 time-varying graphs. Our results introduce these techniques to a large class of real world  
655 applications operating under resource constraints. We provide theoretical guarantees  
656 and numerical validation of the proposed methods in several settings.

657 Future work includes extending these results to non-convex settings with more  
658 sophisticated tasks such as language modelling and speech processing. Moreover, it is  
659 of interest to study stochastic gradient methods as they will reduce local computations.  
660 Finally, an orthogonal direction worth exploring involves security models that account  
661 for adversarial attacks.

662

## REFERENCES

- 663 [1] M. ABADI, A. CHU, I. GOODFELLOW, H. B. McMAHAN, I. MIRONOV, K. TALWAR, AND  
664 L. ZHANG, *Deep learning with differential privacy*, in Proceedings of the 2016 ACM SIGSAC  
665 Conference on Computer and Communications Security, 2016, pp. 308–318.  
666 [2] N. AGARWAL, A. T. SURESH, F. X. X. YU, S. KUMAR, AND B. McMAHAN, *cpsgd: Communication-efficient and differentially-private distributed sgd*, in  
667 Advances in Neural Information Processing Systems 31, S. Bengio, H. Wal-  
668 lach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, eds.,  
669

- 670 Curran Associates, Inc., 2018, pp. 7564–7575, [http://papers.nips.cc/paper/](http://papers.nips.cc/paper/7984-cpsgd-communication-efficient-and-differentially-private-distributed-sgd.pdf)  
671 [7984-cpsgd-communication-efficient-and-differentially-private-distributed-sgd.pdf](http://papers.nips.cc/paper/7984-cpsgd-communication-efficient-and-differentially-private-distributed-sgd.pdf).
- 672 [3] M. ASSRAN, N. LOIZOU, N. BALLAS, AND M. RABBAT, *Stochastic gradient push for distributed*  
673 *deep learning*, arXiv preprint arXiv:1811.10792, (2018).
- 674 [4] A. BELLET, R. GUERRAOU, M. TAZIKI, AND M. TOMMASI, *Personalized and private peer-to-*  
675 *peer machine learning*, arXiv preprint arXiv:1705.08435, (2017).
- 676 [5] S. BOYD, A. GHOSH, B. PRABHAKAR, AND D. SHAH, *Gossip algorithms: design, analysis and*  
677 *applications*, in Proceedings IEEE 24th Annual Joint Conference of the IEEE Computer  
678 and Communications Societies., vol. 3, 2005, pp. 1653–1664 vol. 3.
- 679 [6] K. CAI AND H. ISHII, *Average consensus on general strongly connected digraphs*, *Automatica*,  
680 48 (2012), pp. 2750–2761.
- 681 [7] K. CAI AND H. ISHII, *Average consensus on arbitrary strongly connected digraphs with time-*  
682 *varying topologies*, *IEEE Transactions on Automatic Control*, 59 (2014), pp. 1066–1071.
- 683 [8] T.-H. CHANG, M. HONG, AND X. WANG, *Multi-agent distributed optimization via inexact*  
684 *consensus admm*, *IEEE Transactions on Signal Processing*, 63 (2014), pp. 482–497.
- 685 [9] W.-N. CHEN, P. KAIROUZ, AND A. ÖZGÜR, *Breaking the communication-privacy-accuracy*  
686 *trilemma*, arXiv preprint arXiv:2007.11707, (2020).
- 687 [10] Y. CHEN, A. HASHEMI, AND H. VIKALO, *Communication-efficient algorithms for decentralized*  
688 *optimization over directed graphs*, 2020, <https://arxiv.org/abs/2005.13189>.
- 689 [11] J. CORTÉS, G. E. DULLERUD, S. HAN, J. LE NY, S. MITRA, AND G. J. PAPPAS, *Differential*  
690 *privacy in control and network systems*, in 2016 IEEE 55th Conference on Decision and  
691 Control (CDC), 2016, pp. 4252–4272.
- 692 [12] J. DAVIN AND P. J. MODI, *Impact of problem centralization in distributed constraint optimization*  
693 *algorithms*, in Proceedings of the fourth international joint conference on Autonomous agents  
694 and multiagent systems, 2005, pp. 1057–1063.
- 695 [13] A. D. DOMÍNGUEZ-GARCÍA AND C. N. HADJICOSTIS, *Distributed strategies for average consensus*  
696 *in directed graphs*, in 2011 50th IEEE Conference on Decision and Control and European  
697 Control Conference, IEEE, 2011, pp. 2124–2129.
- 698 [14] J. C. DUCHI, A. AGARWAL, AND M. J. WAINWRIGHT, *Dual averaging for distributed optimiza-*  
699 *tion: Convergence analysis and network scaling*, *IEEE Transactions on Automatic Control*,  
700 57 (2012), pp. 592–606.
- 701 [15] C. DWORK, F. MCSHERRY, K. NISSIM, AND A. SMITH, *Calibrating noise to sensitivity in*  
702 *private data analysis*, Springer, 2006, pp. 265–284.
- 703 [16] C. DWORK, A. ROTH, ET AL., *The algorithmic foundations of differential privacy*, *Foundations*  
704 *and Trends® in Theoretical Computer Science*, 9 (2014), pp. 211–407.
- 705 [17] P. ERDÖS, A. RÉNYI, ET AL., *On random graphs*, *Publicationes mathematicae*, 6 (1959),  
706 pp. 290–297.
- 707 [18] A. HARD, K. RAO, R. MATHEWS, S. RAMASWAMY, F. BEAUFAYS, S. AUGENSTEIN, H. EICHNER,  
708 C. KIDDON, AND D. RAMAGE, *Federated learning for mobile keyboard prediction*, arXiv  
709 preprint arXiv:1811.03604, (2018).
- 710 [19] J. HE, L. CAI, C. ZHAO, P. CHENG, AND X. GUAN, *Privacy-preserving average consensus:*  
711 *Privacy analysis and algorithm design*, *IEEE Transactions on Signal and Information*  
712 *Processing over Networks*, 5 (2019), pp. 127–138.
- 713 [20] A. JADBABAIE, JIE LIN, AND A. S. MORSE, *Coordination of groups of mobile autonomous*  
714 *agents using nearest neighbor rules*, *IEEE Transactions on Automatic Control*, 48 (2003),  
715 pp. 988–1001.
- 716 [21] P. KAIROUZ, H. B. MCMAHAN, B. AVENT, A. BELLET, M. BENNIS, A. N. BHAGOJI,  
717 K. BONAWITZ, Z. CHARLES, G. CORMODE, R. CUMMINGS, ET AL., *Advances and open*  
718 *problems in federated learning*, arXiv preprint arXiv:1912.04977, (2019).
- 719 [22] P. KAIROUZ, S. OH, AND P. VISWANATH, *The composition theorem for differential privacy*, in  
720 International conference on machine learning, PMLR, 2015, pp. 1376–1385.
- 721 [23] S. P. KASIVISWANATHAN, H. K. LEE, K. NISSIM, S. RASKHODNIKOVA, AND A. SMITH, *What*  
722 *can we learn privately?*, *SIAM Journal on Computing*, 40 (2011), pp. 793–826.
- 723 [24] A. KOLOSKOVA, T. LIN, S. U. STICH, AND M. JAGGI, *Decentralized deep learning with arbitrary*  
724 *communication compression*, arXiv preprint arXiv:1907.09356, (2019).
- 725 [25] X. LI, W. YANG, S. WANG, AND Z. ZHANG, *Communication efficient decentralized training*  
726 *with multiple local updates*, arXiv preprint arXiv:1910.09126, (2019).
- 727 [26] Y. LI, L. GUO, AND S. K. PRASAD, *An energy-efficient distributed algorithm for minimum-*  
728 *latency aggregation scheduling in wireless sensor networks*, in 2010 IEEE 30th International  
729 Conference on Distributed Computing Systems, 2010, pp. 827–836.
- 730 [27] X. LIAN, C. ZHANG, H. ZHANG, C.-J. HSIEH, W. ZHANG, AND J. LIU, *Can decentralized algo-*  
731 *rithms outperform centralized algorithms? a case study for decentralized parallel stochastic*

- 732 *gradient descent*, in Advances in Neural Information Processing Systems 30, I. Guyon, U. V.  
 733 Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, eds., Curran  
 734 Associates, Inc., 2017, pp. 5330–5340.
- [28] X. LIAN, C. ZHANG, H. ZHANG, C.-J. HSIEH, W. ZHANG, AND J. LIU, *Can decentralized*  
 736 *algorithms outperform centralized algorithms? a case study for decentralized parallel sto-*  
 737 *chastic gradient descent*, in Advances in Neural Information Processing Systems, 2017,  
 738 pp. 5330–5340.
- [29] H. B. McMAHAN, G. ANDREW, U. ERLINGSSON, S. CHIEN, I. MIRONOV, N. PAPERNOT,  
 740 AND P. KAIROUZ, *A general approach to adding differential privacy to iterative training*  
 741 *procedures*, arXiv preprint arXiv:1812.06210, (2018).
- [30] H. B. McMAHAN, E. MOORE, D. RAMAGE, S. HAMPSON, ET AL., *Communication-efficient*  
 743 *learning of deep networks from decentralized data*, arXiv preprint arXiv:1602.05629, (2016).
- [31] I. MIRONOV, *Rényi differential privacy*, in 2017 IEEE 30th Computer Security Foundations  
 744 Symposium (CSF), IEEE, 2017, pp. 263–275.
- [32] I. MIRONOV, K. TALWAR, AND L. ZHANG, *Rényi differential privacy of the sampled gaussian*  
 746 *mechanism*, arXiv preprint arXiv:1908.10530, (2019).
- [33] A. NARAYANAN AND V. SHMATIKOV, *Robust de-anonymization of large sparse datasets*, in 2008  
 748 IEEE Symposium on Security and Privacy (sp 2008), IEEE, 2008, pp. 111–125.
- [34] A. NEDIĆ, S. LEE, AND M. RAGINSKY, *Decentralized online optimization with global objectives*  
 750 *and local communication*, in 2015 American Control Conference (ACC), IEEE, 2015, pp. 4497–  
 751 4503.
- [35] A. NEDIĆ AND A. OLSHEVSKY, *Distributed optimization over time-varying directed graphs*,  
 753 IEEE Transactions on Automatic Control, 60 (2014), pp. 601–615.
- [36] A. NEDIĆ AND A. OLSHEVSKY, *Stochastic gradient-push for strongly convex functions on time-*  
 755 *varying directed graphs*, IEEE Transactions on Automatic Control, 61 (2016), pp. 3936–3947.
- [37] A. NEDIC, A. OLSHEVSKY, AND W. SHI, *Achieving geometric convergence for distributed*  
 757 *optimization over time-varying graphs*, SIAM Journal on Optimization, 27 (2017), pp. 2597–  
 758 2633.
- [38] A. NEDIC AND A. OZDAGLAR, *Distributed subgradient methods for multi-agent optimization*,  
 760 IEEE Transactions on Automatic Control, 54 (2009), pp. 48–61.
- [39] E. NOZARI, P. TALLAPRAGADA, AND J. CORTÉS, *Differentially private average consensus:*  
 762 *Obstructions, trade-offs, and optimal algorithm design*, Automatica, 81 (2017), pp. 221–231.
- [40] A. PAVERD, A. MARTIN, AND I. BROWN, *Modelling and automatically analysing privacy*  
 764 *properties for honest-but-curious adversaries*.
- [41] S. REDDI, Z. CHARLES, M. ZAHEER, Z. GARRETT, K. RUSH, J. KONEČNÝ, S. KUMAR,  
 766 AND H. B. McMAHAN, *Adaptive federated optimization*, arXiv preprint arXiv:2003.00295,  
 767 (2020).
- [42] W. REN AND R. W. BEARD, *Consensus seeking in multiagent systems under dynamically*  
 769 *changing interaction topologies*, IEEE Transactions on automatic control, 50 (2005), pp. 655–  
 770 661.
- [43] W. REN, H. CHAO, W. BOURGEOUS, N. SORENSSEN, AND Y. CHEN, *Experimental validation*  
 772 *of consensus algorithms for multivehicle cooperative control*, IEEE Transactions on Control  
 773 Systems Technology, 16 (2008), pp. 745–752.
- [44] D. SHAH, *Gossip algorithms*, Now Publishers Inc, 2009.
- [45] S. U. STICH, J.-B. CORDONNIER, AND M. JAGGI, *Sparsified sgd with memory*, in Advances in  
 776 Neural Information Processing Systems, 2018, pp. 4447–4458.
- [46] A. T. SURESH, F. X. YU, S. KUMAR, AND H. B. McMAHAN, *Distributed mean estimation*  
 778 *with limited communication*, in Proceedings of the 34th International Conference on Machine  
 779 Learning-Volume 70, JMLR. org, 2017, pp. 3329–3337.
- [47] H. TANG, S. GAN, C. ZHANG, T. ZHANG, AND J. LIU, *Communication compression*  
 781 *for decentralized training*, in Advances in Neural Information Processing Systems 31,  
 782 S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Gar-  
 783 nett, eds., Curran Associates, Inc., 2018, pp. 7652–7662, [http://papers.nips.cc/paper/](http://papers.nips.cc/paper/7992-communication-compression-for-decentralized-training.pdf)  
 784 [7992-communication-compression-for-decentralized-training.pdf](http://papers.nips.cc/paper/7992-communication-compression-for-decentralized-training.pdf).
- [48] H. TANG, S. GAN, C. ZHANG, T. ZHANG, AND J. LIU, *Communication compression for*  
 786 *decentralized training*, in Advances in Neural Information Processing Systems, 2018, pp. 7652–  
 787 7662.
- [49] J. WANG AND G. JOSHI, *Cooperative sgd: A unified framework for the design and analysis of*  
 789 *communication-efficient sgd algorithms*, arXiv preprint arXiv:1808.07576, (2018).
- [50] J. WANG, A. K. SAHU, Z. YANG, G. JOSHI, AND S. KAR, *Matcha: Speeding up decentralized*  
 791 *sgd via matching decomposition sampling*, arXiv preprint arXiv:1905.09435, (2019).
- [51] Y. WANG, *Privacy-preserving average consensus via state decomposition*, IEEE Transactions

- 794           on Automatic Control, 64 (2019), pp. 4711–4716.
- 795 [52] E. WEI AND A. OZDAGLAR, *Distributed alternating direction method of multipliers*, in 2012  
796 IEEE 51st IEEE Conference on Decision and Control (CDC), IEEE, 2012, pp. 5445–5450.
- 797 [53] WEI REN AND R. W. BEARD, *Consensus seeking in multiagent systems under dynamically*  
798 *changing interaction topologies*, IEEE Transactions on Automatic Control, 50 (2005),  
799 pp. 655–661.
- 800 [54] X. WU, F. LI, A. KUMAR, K. CHAUDHURI, S. JHA, AND J. NAUGHTON, *Bolt-on differential*  
801 *privacy for scalable stochastic gradient descent-based analytics*, ACM, 2017, pp. 1307–1322.
- 802 [55] C. XI, Q. WU, AND U. A. KHAN, *On the distributed optimization over directed networks*,  
803 Neurocomputing, 267 (2017), pp. 508–515.