

---

# Central Server Free Federated Learning over Single-sided Trust Social Networks

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Federated learning has become increasingly important for modern machine learning,  
2 especially for data privacy-sensitive scenarios. Existing federated learning primarily  
3 adopts the central server-based network topology. However, in many social network  
4 scenarios, centralized federated learning is not applicable. For instance, a central  
5 agent or server connecting all users may not exist, or the communication cost to  
6 the central server is not affordable. In this paper, we consider a generic setting: 1)  
7 the central server may not exist, and 2) the social network is unidirectional or of  
8 single-sided trust (i.e., user A trusts user B but user B may not trust user A). We  
9 propose a central server free federated learning algorithm, named Online Push-Sum  
10 (OPS) method, to handle this challenging but generic scenario. A rigorous regret  
11 analysis is also provided, which shows interesting results on how users can benefit  
12 from communication with trusted users in the federated learning scenario. This  
13 work builds upon the fundamental algorithm framework and theoretical guarantees  
14 for federated learning in the generic social network scenario.

## 15 1 Introduction

16 Federated learning has been well recognized as a framework able to protect data privacy [1, 2, 3].  
17 State-of-the-art federated learning adopts the centralized network architecture where a centralized  
18 node collects the gradients sent from child agents to update the global model. Despite its simplicity,  
19 the centralized method suffers from communication and computational bottlenecks in the central node,  
20 especially for federated learning, where a large number of clients are usually involved. Moreover,  
21 to prevent reverse engineering of the user’s identity, a certain amount of noise must be added to the  
22 gradient to protect user privacy, which partially sacrifices the efficiency and the accuracy [4].

23 To further protect the data privacy and avoid the communication bottleneck, the decentralized  
24 architecture has been recently proposed [5, 6], where the centralized node has been removed, and  
25 each node only communicates with its neighbors (with mutual trust) by exchanging their local models.  
26 Exchanging local models is usually favored to the data privacy protection over sending private  
27 gradients because the local model is the aggregation or mixture of a large sum of data while the  
28 local gradient directly reflects only one or a batch of private data samples. Although advantages of  
29 decentralized architecture have been well recognized over the state-of-the-art method (its centralized  
30 counterpart), it usually can only be run on the network with *mutual trusts* (“trust” means “would like  
31 to send information to”). That is, two nodes (or users) can exchange their local models only if they  
32 trust each other reciprocally (e.g., node A may trust node B, but if node B does not trust node A,  
33 they cannot communicate). Given a social network, one can only use the edges with mutual trust  
34 to run decentralized federated learning algorithms. Two immediate drawbacks are: (1) If all mutual  
35 trust edges do not form a connected network, the federated learning does not apply; (2) Removing  
36 all single-sided edges from the communication network could significantly reduce the efficiency of

37 communication. These drawbacks lead to the question: *How do we effectively utilize the single-sided*  
 38 *trust edges under a decentralized federated learning framework?*

39 In this paper, we consider the social network scenario, where the centralized network is unavailable  
 40 (e.g., there does not exist a central node that can build up the connection with all users, or the  
 41 centralized communication cost is not affordable). We make a minimal assumption on the social  
 42 network: The data may come in a streaming fashion on each user node as the federated learning  
 43 algorithm runs; the trust between users may be single-sided, where user A trusts user B, but user B  
 44 may not trust user A.

45 For the setting mentioned above, we develop a decentralized learning algorithm called online push-  
 46 sum (OPS) which possesses the following features:

- 47 • Our algorithm removes some constraints imposed by typical decentralized methods, making  
 48 it more flexible in allowing arbitrary network topology. Each node only needs to know its  
 49 out neighbors rather than the global topology.
- 50 • Only models rather than local gradients are exchanged among clients in our algorithm. This  
 51 scheme can reduce the risk of exposing clients' data privacy [7].
- 52 • We provide the rigorous regret analysis for the proposed algorithm and specifically dis-  
 53 tinguish two components in the online loss function: the adversary component and the  
 54 stochastic component, which can model clients' private data and internal connections be-  
 55 tween clients, respectively.

56 **Notation** We adopt the following notation in this paper:

- 57 • For random variable  $\xi_t^{(i)}$  subject to distribution  $D_t^{(i)}$ , we use  $\Xi_{n,T}$  and  $\mathcal{D}_{n,T}$  to denote the  
 58 set of random variables and distributions, respectively:

$$\Xi_{n,T} = \{\xi_t^{(i)}\}_{1 \leq i \leq n, 1 \leq t \leq T}, \quad \mathcal{D}_{n,T} = \{D_t^{(i)}\}_{1 \leq i \leq n, 1 \leq t \leq T}.$$

59 Notation  $\Xi_{n,T} \sim \mathcal{D}_{n,T}$  implies  $\xi_t^{(i)} \sim D_t^{(i)}$  for any  $i \in [n]$  and  $t \in [T]$ .

- 60 • For a decentralized network with  $n$  nodes, we use  $\mathbf{W} \in \mathbb{R}^{n \times n}$  to present the confusion  
 61 matrix, where  $W_{ij} \geq 0$  is the weight that node  $i$  sends to node  $j$  ( $i, j \in [n]$ ).  $\mathcal{N}_i^{\text{out}} = \{j \in$   
 62  $[n] : W_{ij} > 0\}$  and  $\mathcal{N}_i^{\text{in}} = \{k \in [n] : W_{ki} > 0\}$  are also used for denoting the sets of in  
 63 neighbors of and out neighbors of node  $i$  respectively.
- 64 • Norm  $\|\cdot\|$  denotes the  $\ell_2$  norm  $\|\cdot\|_2$  by default.

## 65 2 Problem Setting

66 In this paper, we consider federated learning with  $n$  clients (a.k.a., nodes). Each client can be either  
 67 an edge server or some other kind of computing device such as a smartphone, which has local private  
 68 data and the local machine learning model  $\mathbf{x}_i$  stored on it. We assume the topological structure of  
 69 the network of these  $n$  nodes can be represented by a directed graph  $\mathcal{G} = (\text{nodes} : [n], \text{edges} : E)$   
 70 with vertex set  $[n] = \{1, 2, \dots, n\}$  and edge set  $E \subset [n] \times [n]$ . If there exists an edge  $(u, v) \in E$ , it  
 71 means node  $u$  and node  $v$  have network connection and  $u$  can directly send messages to  $v$ .

72 Let  $\mathbf{x}_t^{(i)}$  denote the local model on the  $i$ -th node at iteration  $t$ . In each iteration, node  $i$  receives a new  
 73 sample and computes a prediction for this new sample according to the current model  $\mathbf{x}_t^{(i)}$  (e.g., it  
 74 may recommend some items to the user in the online recommendation system). After that, a loss  
 75 function,  $f_{i,t}(\cdot)$  associated with that new sample is received by node  $i$ . The typical goal of online  
 76 learning is to minimize the *regret*, which is defined as the difference between the summation of the  
 77 losses incurred by the nodes' prediction and the corresponding loss of the global optimal model  $\mathbf{x}^*$ :

$$\tilde{\mathcal{R}}_T := \sum_{t=1}^T \sum_{i=1}^n \left( f_{i,t}(\mathbf{x}_t^{(i)}) - f_{i,t}(\mathbf{x}^*) \right),$$

78 where  $\mathbf{x}^* = \arg \min_{\mathbf{x}} \sum_{t=1}^T \sum_{i=1}^n f_{i,t}(\mathbf{x})$  is the optimal solution.

79 However, here we consider a more general online setting: the loss function of the  $i$ -th node at iteration  
80  $t$  is  $f_{i,t}(\cdot; \xi_{i,t})$ , which is additionally parametrized by a random variable  $\xi_{i,t}$ . This  $\xi_{i,t}$  is drawn  
81 from the distribution  $D_{i,t}$ , and is mutually independent in terms of  $i$  and  $t$ , and we call this part as  
82 the *stochastic* component of loss function  $f_{i,t}(\cdot; \xi_{i,t})$ . The stochastic component can be utilized to  
83 characterize the internal randomness of nodes' data, and the potential connection among different  
84 nodes. For example, music preference may be impacted by popular trends on the Internet, which can  
85 be formulated by our model by letting  $D_{i,t} \equiv D_t$  for all  $i \in [n]$  with some time-varying distribution  
86  $D_t$ . On the other hand, function  $f_{i,t}(\cdot; \cdot)$  is the *adversarial* component of the loss, which may include,  
87 for example, user's profile, location, etc. Therefore, the objective regret naturally becomes the  
88 expectation of all the past losses:

$$\mathcal{R}_T := \mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \left\{ \sum_{t=1}^T \sum_{i=1}^n \left( f_{i,t}(\mathbf{x}_t^{(i)}; \xi_t^{(i)}) - f_{i,t}(\mathbf{x}^*; \xi_t^{(i)}) \right) \right\} \quad (1)$$

89 with  $\mathbf{x}^* = \arg \min_{\mathbf{x}} \mathbb{E}_{\Xi_{n,T} \sim \mathcal{D}_{n,T}} \sum_{t=1}^T \sum_{i=1}^n f_{i,t}(\mathbf{x}; \xi_t^{(i)})$ .

90 One benefit of the above formulation is that it partially resolves the non-I.I.D. issue in federated  
91 learning. A fundamental assumption in many traditional distributed machine learning methods is  
92 that the data samples stored on all nodes are I.I.D., which fails to hold for federated learning since  
93 the data on each user's device is highly correlated to that user's preferences and habits. However,  
94 our formulation does not require the I.I.D. assumption to hold for the adversarial component at all.  
95 Even though the random samples for the stochastic component still need to be independent, they are  
96 allowed to be drawn from different distributions.

97 Finally, one should note that online optimization also includes stochastic optimization (i.e., data  
98 samples are drawn from a fixed distribution) and offline optimization (i.e., data are already collected  
99 before optimization begins) as its typical cases [8]. Hence, our setting covers a wide range of  
100 applications.

### 101 3 Online Push-Sum Algorithm

102 In this section, we define the construction of the confusion matrix and introduce the proposed  
103 algorithm.

#### 104 3.1 Construction of Confusion Matrix

105 One important parameter of the algorithm is the confusion matrix  $\mathbf{W}$ .  $\mathbf{W}$  is a matrix depending on  
106 the network topology  $\mathcal{G}$ , which means  $W_{ij} = 0$  if there is no directed edge  $(i, j)$  in  $\mathcal{G}$ . If the value of  
107  $W_{ij}$  is large, the node  $i$  will have a stronger impact on node  $j$ . However,  $\mathbf{W}$  still allows flexibility  
108 where users can specify their weights associated with existing edges, meaning that even if there is a  
109 physical connection between two nodes, the nodes can decide against using the channel. For example,  
110 even if  $(i, j) \in E$ , user still can set  $W_{ij} = 0$  if user  $i$  thinks node  $j$  is not trustworthy and therefore  
111 chooses to exclude the channel from  $i$  to  $j$ .

112 Of course, there are still some constraints over  $\mathbf{W}$ .  $\mathbf{W}$  must be a row stochastic matrix (i.e., each  
113 entry in  $\mathbf{W}$  is non-negative, and the summation of each row is 1). This assumption is different from  
114 the one in classical decentralized distributed optimization, which typically assumes  $\mathbf{W}$  is symmetric  
115 and doubly stochastic (e.g., [9]) (i.e., the summations of both rows and columns are all 1). Such a  
116 requirement is quite restrictive, because not all networks admit a doubly stochastic matrix ([10]), and  
117 relinquishing double stochasticity can introduce bias in optimization [11, 12]. As a comparison, our  
118 assumption that  $\mathbf{W}$  is row stochastic will avoid such concerns since any non-negative matrix with at  
119 least one positive entry on each row (which is already implied by the connectivity of the graph) can  
120 be easily normalized into row stochastic. The relaxation of this assumption is crucial for federated  
121 learning, considering that the federated learning system usually involves complex network topology  
122 due to its large number of clients. Moreover, since each node only needs to make sure the summation  
123 of its out-weights is 1, there is no need for it to be aware of the global network topology, which  
124 significantly benefits the implementation of the federated learning system. Meanwhile, requiring  
125  $\mathbf{W}$  to be symmetric rules out the possibility of using asymmetric network topology and adopting  
126 sing-sided trust, while our method does not have such restriction.

---

**Algorithm 1** Online Push-Sum (OPS) Algorithm
 

---

**Require:** Learning rate  $\gamma$ , number of iterations  $T$ , and the confusion matrix  $\mathbf{W}$ .

|   |   |
|---|---|
| <p>1: Initialize <math>\mathbf{x}_0^{(i)} = \mathbf{z}_0^{(i)} = \mathbf{0}</math>, <math>\omega_0^{(i)} = 1</math> for all <math>i \in [n]</math></p> <p>2: <b>for</b> <math>t = 0, 1, \dots, T - 1</math> <b>do</b></p> <p>3:     // For all users (say the <math>i</math>-th node <math>i \in [n]</math>)</p> <p>4:     Apply local model <math>\mathbf{x}_t^{(i)}</math> and suffer loss <math>f_{i,t}(\mathbf{x}_t^{(i)}; \boldsymbol{\xi}_t^{(i)})</math></p> <p>5:     Locally computes the intermedia variable <math>\mathbf{z}_{t+\frac{1}{2}}^{(i)} = \mathbf{z}_t^{(i)} - \gamma \nabla f_{i,t}(\mathbf{x}_t^{(i)}; \boldsymbol{\xi}_t^{(i)})</math></p> <p>6:     Send <math>(W_{ij}\mathbf{z}_{t+\frac{1}{2}}^{(i)}, W_{ij}\omega_t^{(i)})</math> to all <math>j \in \mathcal{N}_i^{\text{out}}</math></p> | <p>7:     Update</p> $\mathbf{z}_{t+1}^{(i)} = \sum_{k \in \mathcal{N}_i^{\text{in}}} W_{ki} \mathbf{z}_{t+\frac{1}{2}}^{(k)}$ $\omega_{t+1}^{(i)} = \sum_{k \in \mathcal{N}_i^{\text{in}}} W_{ki} \omega_t^{(k)}$ $\mathbf{x}_{t+1}^{(i)} = \frac{\mathbf{z}_{t+1}^{(i)}}{\omega_{t+1}^{(i)}}$ <p>8: <b>end for</b></p> <p>9: <b>return</b> <math>\mathbf{x}_T^{(i)}</math> to node <math>i</math></p> |
|---|---|

---

### 127 3.2 Algorithm Description

128 The proposed online push-sum algorithm is presented in Algorithm 1. The algorithm design mainly  
 129 follows the pattern of push-sum algorithm [13], but here we further generalize it into the online  
 130 setting.

131 The algorithm mainly consists of three steps:

- 132     1. Local update: each client  $i$  applies the current local model  $\mathbf{x}_t^{(i)}$  to obtain the loss function,  
 133         based on which an intermediate local model  $\mathbf{z}_{t+\frac{1}{2}}^{(i)}$  is computed;
- 134     2. Push: the weighted variable  $W_{ij}\mathbf{z}_{t+\frac{1}{2}}^{(i)}$  is sent to  $j$  for all its out neighbors  $j$ ;
- 135     3. Sum: all the received  $W_{ji}\mathbf{z}_{t+\frac{1}{2}}^{(j)}$  is summed and normalized to obtain the new model  $\mathbf{x}_{t+1}^{(i)}$ .

136 It should be noted an auxiliary variables  $\mathbf{z}_{t+\frac{1}{2}}^{(i)}$  and  $\mathbf{z}_{t+1}^{(i)}$  are used in the algorithm. Actually, they  
 137 are used in the algorithm to clarify the description but may be easily removed in the practical  
 138 implementation. Besides, another variable  $\omega_{t+1}^{(i)}$  is also introduced, which is the normalizing factor of  
 139  $\mathbf{z}_{t+1}^{(i)}$ .  $\omega_{t+1}^{(i)}$  plays an important role in the push-sum algorithm, since  $\mathbf{W}$  is not doubly stochastic in  
 140 our setting, and it is possible that the total weight  $i$  receives does not equal to 1. The introduction  
 141 of the normalizing factor  $\omega_t^{(i)}$  helps the algorithm avoid issues brought by that  $\mathbf{W}$  is not doubly  
 142 stochastic. Furthermore, when  $\mathbf{W}$  becomes doubly stochastic, it can be easily verified that  $\omega_t^{(i)} \equiv 1$   
 143 and  $\mathbf{x}_t^{(i)} \equiv \mathbf{z}_t^{(i)}$  for any  $i$  and  $t$ , then Algorithm 1 reduces to the distributed online gradient method  
 144 proposed by [14].

### 145 3.3 Regret Analysis

146 In this subsection, we provide regret bound analysis of OPS algorithm. Due to the limitation of space,  
 147 the detail proof is deferred to the supplementary material. For convenience, we first denote

$$F_{i,t}(\mathbf{x}) := \mathbb{E}_{\boldsymbol{\xi}_{i,t} \sim D_{i,t}} f_{i,t}(\mathbf{x}; \boldsymbol{\xi}_{i,t}).$$

148 To carry out the analysis, the following assumptions are required:

149 **Assumption 1.** We make the following assumptions throughout this paper: (1) The topological graph  
 150  $\mathcal{G}$  is strongly connected;  $\mathbf{W}$  is row stochastic; (2) For any  $i \in [n]$  and  $t \in [T]$ , the loss function  
 151  $f_{i,t}(\mathbf{x}; \boldsymbol{\xi}_{i,t})$  is convex in  $\mathbf{x}$ ; (3) The problem domain is bounded such that for any two vectors  $\mathbf{x}$  and  
 152  $\mathbf{y}$  we always have  $\|\mathbf{x} - \mathbf{y}\|^2 \leq R$ ; (4) The norm of the expected gradient  $\nabla F_{i,t}(\cdot)$  is bounded, i.e.,

153 there exist constant  $G > 0$  such that  $\|\nabla F_{i,t}(\mathbf{x})\|^2 \leq G^2$  for any  $i, t$  and  $\mathbf{x}$ ; (5) The gradient variance  
 154 is also bounded by  $\sigma^2$ , namely,

$$\mathbb{E}_{\xi_{i,t} \sim D_{i,t}} \|\nabla f_{i,t}(\mathbf{x}; \xi_{i,t}) - \nabla F_{i,t}(\mathbf{x})\|^2 \leq \sigma^2.$$

155 Equipped with these assumptions, now we are ready to present the convergence result:

156 **Theorem 2.** *If we set*

$$\gamma = \frac{\sqrt{n}R}{\sigma\sqrt{1+nC_2} + G\sqrt{nC_1T}}, \quad (2)$$

157 *the regret of OPS can be bounded by:*

$$\mathcal{R}_T \leq \mathcal{O}\left(nGR\sqrt{T} + \sigma R\left(1 + \sqrt{nC_2}\right)\sqrt{nT}\right), \quad (3)$$

158 *where  $C_1$  and  $C_2$  are two constants defined in the appendix.*

159 Note that when  $n = 1$  and  $\sigma = 0$ , where the problem setting just reduces to normal online  
 160 optimization, the implied regret bound  $\mathcal{O}(GR\sqrt{T})$  exactly matches the lower bound of online  
 161 optimization [15]. Moreover, our result also matches the convergence rate of centralized online  
 162 learning where  $q = 0$  for fully connected networks. Hence, we can conclude that the OPS algorithm  
 163 has optimal dependence on  $T$ .

164 Moreover, we also prove that the difference of the model  $\mathbf{x}_t^{(i)}$  on each worker could be bounded using  
 165 the following theorem:

166 **Theorem 3.** *If we set  $\gamma$  as (2), the difference of the model  $\mathbf{x}_t^{(i)}$  on each worker admits a faster  
 167 convergence rate than regret:*

$$\frac{1}{T} \sum_i^n \sum_{t=0}^T \left\| \mathbf{x}_{t+1}^{(i)} - \bar{\mathbf{z}}_{t+1} \right\|^2 \leq \mathcal{O}\left(\frac{nGR + nR\sigma}{T}\right).$$

168 Hence, the models on all clients' devices will finally converge to the same one with rate  $\mathcal{O}(1/T)$ .

## 169 4 Experiments

170 We compare the performance of our proposed Online Push-Sum (OPS) method with that of Decen-  
 171 tralized Online Gradient method (DOL) and Centralized Online Gradient method (COL), and then  
 172 evaluate the effectiveness of OPS in different network size and network topology density settings.

### 173 4.1 Implementation and Settings

174 We consider online logistic regression with squared  $\ell_2$  norm regularization:  $f_{i,t}(\mathbf{x}; \xi_{i,t}) =$   
 175  $\log(1 + \exp(-\mathbf{y}_{i,t} \mathbf{A}_{i,t}^\top \mathbf{x})) + \frac{\lambda}{2} \|\mathbf{x}\|^2$ , where regularization coefficient  $\lambda$  is set to  $10^{-4}$ .  $\xi_{i,t}$  is  
 176 the stochastic component of the function  $f_{i,t}$  introduced in Section § 2, which is encoded in the  
 177 random data sample  $(\mathbf{A}_{i,t}, \mathbf{y}_{i,t})$ . We evaluate the learning performance by measuring the average  
 178 loss  $\frac{1}{nT} \mathbb{E}_{\Xi_{n,T}} \sum_{i=1}^n \sum_{t=1}^T f_{i,t}(\mathbf{x}_{i,t}; \xi_{i,t})$ , instead of using the dynamic regret (1) directly, since the  
 179 optimal reference point  $x^*$  is the same for all the methods. The learning rate  $\gamma$  in Algorithm 1 is  
 180 tuned to be optimal for each dataset separately.

181 **Dataset.** Experiments were run on two real-world public datasets: *SUSY*<sup>1</sup> and *Room-Occupancy*<sup>2</sup>.  
 182 *SUSY* and *Room-Occupancy* are both large-scale binary classification datasets, containing 5,000,000  
 183 and 20,566 samples, respectively. Each dataset is split into two subsets: the stochastic data and the  
 184 adversarial data. The stochastic data is generated by allocating a fraction of samples (e.g., 50% of the  
 185 whole dataset) to nodes randomly and uniformly. The adversarial data is generated by conducting  
 186 on the remaining dataset to produce  $n$  clusters and then allocating every cluster to a node. As we  
 187 analyzed previously, only the scattered stochastic data can boost the model performance by intra-node  
 188 communication. For each node, this pre-acquired data is transformed into streaming data to simulate  
 189 online learning.

<sup>1</sup><https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/binary.html#SUSY>

<sup>2</sup><https://archive.ics.uci.edu/ml/datasets/Occupancy+Detection+>

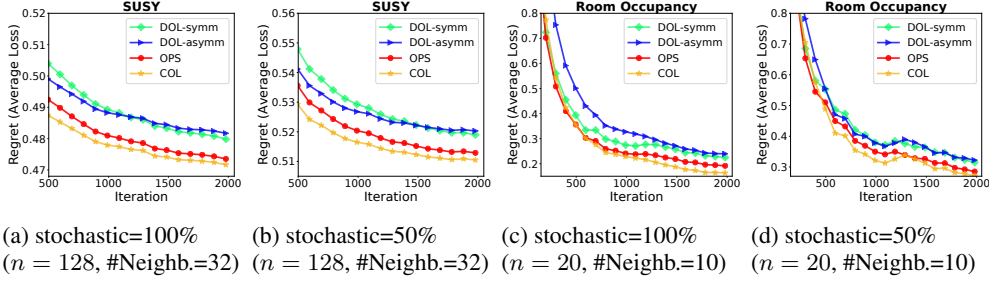


Figure 1: Comparison among OPS, DOL (Decentralized Online Learning) and COL (Centralized Online Learning)

## 190 4.2 Comparison with DOL and COL

191 To compare OPS with DOL and COL, a network size with 128 nodes and 20 nodes are selected for  
 192 SUSY and Room-Occupancy, respectively. For COL, its confusion matrix  $\mathbf{W}$  is fully-connected  
 193 (doubly stochastic matrix). For DOL and OPS, they are run with the same network topology and  
 194 the same row stochastic matrix (asymmetric confusion matrix) to maintain a fair comparison. Such  
 195 asymmetric matrix is constructed by setting each node’s number of neighbors as a random value  
 196 which is smaller than a fixed upper bound and also ensures the strong connectivity of the whole  
 197 network (this upper-bound neighbor number is set to 32 for the SUSY dataset, while 10 is set for  
 198 the Room-Occupancy dataset). Since DOL typically requires the network to be the symmetric and  
 199 doubly stochastic confusion matrix, DOL is run in two settings for comparison. In the first setting, in  
 200 order to meet the assumption of the symmetry and doubly stochasticity, all unidirectional connections  
 201 are removed in the confusion matrix so that the row stochastic confusion matrix degenerates into a  
 202 doubly stochastic matrix. This setting is labeled as *DOL-Symm* in Figure 1. In another setting, DOL  
 203 is forced to run on the asymmetric network where each node naively aggregates its received models  
 204 without considering whether its sending weights are equal to its receiving weights. *DOL-Asymm*  
 205 is used to label this setting in Figure 1.

206 As illustrated in Figure 1, in both two datasets, OPS outperforms *DOL-Symm* in the row stochastic  
 207 confusion matrix. This demonstrates that incorporating unidirectional communication can help to  
 208 boost the model performance. In other words, OPS gains better performance in the single-sided trust  
 209 network under the setting of federated learning. OPS also works better than *DOL-Asymm*. Although  
 210 *DOL-Asymm* utilizes additional unidirectional connections, in some cases its performance is even  
 211 worse than *DOL-Symm* (e.g., Figure 1a). This phenomenon is most likely attributed to its simple  
 212 aggregation pattern, which causes decreased performance in *DOL-Asymm* when removing the doubly  
 213 stochastic matrix assumption. These two observations confirm the effectiveness of OPS in a row  
 214 stochastic confusion matrix, which is consistent with our theoretical analysis.

215 Comparing Figure 1c and Figure 1d, we also observe that when increasing the ratio of the stochastic  
 216 component, the average loss (regret) becomes smaller. It is reasonable that OPS achieves slightly  
 217 worse performance than COL because OPS works in a sparsely connected network where information  
 218 exchanging is much less than COL. We use the COL as the baseline in all experiments.

## 219 4.3 Other Experiments

220 Due to the limitation of space, extra experiments are deferred to the appendix.

## 221 5 Conclusions

222 Decentralized federated learning with single-sided trust is a promising framework for solving a wide  
 223 range of problems. In this paper, the online push-sum algorithm is developed for this setting, which  
 224 is able to handle complex network topology and is proven to have an optimal convergence rate. The  
 225 regret-based online problem formulation also extends its applications. We tested the proposed OPS  
 226 algorithm in various experiments, which have empirically justified its efficiency.

227 **References**

- 228 [1] Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh,  
229 and Dave Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv*  
230 *preprint arXiv:1610.05492*, 2016.
- 231 [2] Virginia Smith, Chao-Kai Chiang, Maziar Sanjabi, and Ameet S Talwalkar. Federated multi-task  
232 learning. In *Advances in Neural Information Processing Systems*, pages 4424–4434, 2017.
- 233 [3] Qiang Yang, Yang Liu, Tianjian Chen, and Yongxin Tong. Federated machine learning: Concept  
234 and applications. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 10(2):12,  
235 2019.
- 236 [4] Reza Shokri and Vitaly Shmatikov. Privacy-preserving deep learning. In *Proceedings of the*  
237 *22nd ACM SIGSAC conference on computer and communications security*, pages 1310–1321.  
238 ACM, 2015.
- 239 [5] Paul Vanhaesebrouck, Aurélien Bellet, and Marc Tommasi. Decentralized collaborative learning  
240 of personalized models over networks. In *International Conference on Artificial Intelligence*  
241 *and Statistics (AISTATS)*, 2017.
- 242 [6] Aurélien Bellet, Rachid Guerraoui, Mahsa Taziki, and Marc Tommasi. Personalized and private  
243 peer-to-peer machine learning. In *International Conference on Artificial Intelligence and*  
244 *Statistics*, pages 473–481, 2018.
- 245 [7] Yoshinori Aono, Takuya Hayashi, Lihua Wang, Shiho Moriai, et al. Privacy-preserving deep  
246 learning via additively homomorphic encryption. *IEEE Transactions on Information Forensics*  
247 *and Security*, 13(5):1333–1345, 2017.
- 248 [8] Shai Shalev-Shwartz et al. Online learning and online convex optimization. *Foundations and*  
249 *Trends® in Machine Learning*, 4(2):107–194, 2012.
- 250 [9] John C Duchi, Alekh Agarwal, and Martin J Wainwright. Dual averaging for distributed  
251 optimization: Convergence analysis and network scaling. *IEEE Transactions on Automatic*  
252 *control*, 57(3):592–606, 2011.
- 253 [10] Bahman Ghahesifard and Jorge Cortés. When does a digraph admit a doubly stochastic adjacency  
254 matrix? In *Proceedings of the 2010 American Control Conference*, pages 2440–2445. IEEE,  
255 2010.
- 256 [11] S Sundhar Ram, Angelia Nedić, and Venugopal V Veeravalli. Distributed stochastic subgradient  
257 projection algorithms for convex optimization. *Journal of optimization theory and applications*,  
258 147(3):516–545, 2010.
- 259 [12] Konstantinos I Tsianos and Michael G Rabbat. Distributed dual averaging for convex opti-  
260 mization under communication delays. In *2012 American Control Conference (ACC)*, pages  
261 1067–1072. IEEE, 2012.
- 262 [13] Konstantinos I Tsianos, Sean Lawlor, and Michael G Rabbat. Push-sum distributed dual  
263 averaging for convex optimization. In *2012 IEEE 51st IEEE Conference on Decision and*  
264 *Control (CDC)*, pages 5453–5458. IEEE, 2012.
- 265 [14] Yawei Zhao, Chen Yu, Peilin Zhao, and Ji Liu. Decentralized online learning: Take benefits from  
266 others’ data without sharing your own to track global trend. *arXiv preprint arXiv:1901.10593*,  
267 2019.
- 268 [15] Elad Hazan et al. Introduction to online convex optimization. *Foundations and Trends® in*  
269 *Optimization*, 2(3-4):157–325, 2016.
- 270 [16] H. Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y  
271 Arcas. Communication-Efficient Learning of Deep Networks from Decentralized Data.  
272 *arXiv:1602.05629 [cs]*, February 2016. arXiv: 1602.05629.

- 273 [17] Jakub Konečný, Brendan McMahan, and Daniel Ramage. Federated Optimization: Distributed  
274 Optimization Beyond the Datacenter. *arXiv:1511.03575 [cs, math]*, November 2015. arXiv:  
275 1511.03575.
- 276 [18] Jakub Konečný, H. Brendan McMahan, Felix X. Yu, Peter Richtárik, Ananda Theertha Suresh,  
277 and Dave Bacon. Federated Learning: Strategies for Improving Communication Efficiency.  
278 *arXiv:1610.05492 [cs]*, October 2016. arXiv: 1610.05492.
- 279 [19] Brendan McMahan and Daniel Ramage. Google AI Blog: Federated Learning: Collaborative  
280 Machine Learning without Centralized Training Data, April 2017.
- 281 [20] Virginia Smith, Chao-Kai Chiang, Maziar Sanjabi, and Ameet S. Talwalkar. Federated multi-task  
282 learning. In *Advances in Neural Information Processing Systems*, pages 4424–4434, 2017.
- 283 [21] Sebastian Caldas, Virginia Smith, and Ameet Talwalkar. Federated Kernelized Multi-Task  
284 Learning. *The Conference on Systems and Machine Learning*, page 3, 2018.
- 285 [22] Shai Shalev-Shwartz and Tong Zhang. Stochastic dual coordinate ascent methods for regularized  
286 loss minimization. *Journal of Machine Learning Research*, 14(Feb):567–599, 2013.
- 287 [23] Tianbao Yang. Trading computation for communication: Distributed stochastic dual coordinate  
288 ascent. In *Advances in Neural Information Processing Systems*, pages 629–637, 2013.
- 289 [24] Tianbao Yang, Shenghuo Zhu, Rong Jin, and Yuanqing Lin. Analysis of distributed stochastic  
290 dual coordinate ascent. *arXiv preprint arXiv:1312.1031*, 2013.
- 291 [25] Martin Jaggi, Virginia Smith, Martin Takáč, Jonathan Terhorst, Sanjay Krishnan, Thomas  
292 Hofmann, and Michael I. Jordan. Communication-efficient distributed dual coordinate ascent.  
293 In *Advances in neural information processing systems*, pages 3068–3076, 2014.
- 294 [26] Chenxin Ma, Virginia Smith, Martin Jaggi, Michael I. Jordan, Peter Richtárik, and Martin  
295 Takáč. Adding vs. Averaging in Distributed Primal-Dual Optimization. *arXiv:1502.03508 [cs]*,  
296 February 2015. arXiv: 1502.03508.
- 297 [27] Virginia Smith, Simone Forte, Chenxin Ma, Martin Takac, Michael I. Jordan, and Martin  
298 Jaggi. CoCoA: A General Framework for Communication-Efficient Distributed Optimization.  
299 *arXiv:1611.02189 [cs]*, November 2016. arXiv: 1611.02189.
- 300 [28] Sebastian U. Stich. Local SGD Converges Fast and Communicates Little. September 2018.
- 301 [29] Jianyu Wang and Gauri Joshi. Cooperative SGD: A unified Framework for the Design and  
302 Analysis of Communication-Efficient SGD Algorithms. August 2018.
- 303 [30] Hao Yu, Sen Yang, and Shenghuo Zhu. Parallel Restarted SGD with Faster Convergence  
304 and Less Communication: Demystifying Why Model Averaging Works for Deep Learning.  
305 *arXiv:1807.06629 [cs, math]*, July 2018. arXiv: 1807.06629.
- 306 [31] Tao Lin, Sebastian U. Stich, and Martin Jaggi. Don’t Use Large Mini-Batches, Use Local SGD.  
307 *arXiv:1808.07217 [cs, stat]*, August 2018. arXiv: 1808.07217.
- 308 [32] Xiangru Lian, Ce Zhang, Huan Zhang, Cho-Jui Hsieh, Wei Zhang, and Ji Liu. Can decentralized  
309 algorithms outperform centralized algorithms? a case study for decentralized parallel stochastic  
310 gradient descent. In *Advances in Neural Information Processing Systems*, pages 5330–5340,  
311 2017.
- 312 [33] Lie He, An Bian, and Martin Jaggi. Cola: Decentralized linear learning. In *Advances in Neural  
313 Information Processing Systems*, pages 4541–4551, 2018.
- 314 [34] Angelia Nedić and Alex Olshevsky. Distributed optimization over time-varying directed graphs.  
315 *IEEE Transactions on Automatic Control*, 60(3):601–615, 2015.
- 316 [35] Xiangru Lian, Wei Zhang, Ce Zhang, and Ji Liu. Asynchronous decentralized parallel stochastic  
317 gradient descent. In *International Conference on Machine Learning*, 2018.



- 318 [36] Tianyu Wu, Kun Yuan, Qing Ling, Wotao Yin, and Ali H. Sayed. Decentralized consensus  
319 optimization with asynchrony and delays. *IEEE Transactions on Signal and Information*  
320 *Processing over Networks*, PP:1–1, 04 2017.
- 321 [37] Zebang Shen, Aryan Mokhtari, Tengfei Zhou, Peilin Zhao, and Hui Qian. Towards more  
322 efficient stochastic decentralized learning: Faster convergence and sparse communication. In  
323 Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on*  
324 *Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 4624–4633,  
325 Stockholmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR.
- 326 [38] Youjie Li, Mingchao Yu, Songze Li, Salman Avestimehr, Nam Sung Kim, and Alexander  
327 Schwing. Pipe-sgd: A decentralized pipelined sgd framework for distributed deep net training.  
328 In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors,  
329 *Advances in Neural Information Processing Systems 31*, pages 8056–8067. Curran Associates,  
330 Inc., 2018.
- 331 [39] Michael Kamp, Mario Boley, Daniel Keren, Assaf Schuster, and Izchak Sharfman.  
332 Communication-efficient distributed online prediction by dynamic model synchronization.  
333 In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*,  
334 pages 623–639. Springer, 2014.
- 335 [40] Shahin Shahrapour and Ali Jadbabaie. Distributed online optimization in dynamic envi-  
336 ronments using mirror descent. *IEEE Transactions on Automatic Control*, 63(3):714–725,  
337 2017.
- 338 [41] Soomin Lee, Angelia Nedić, and Maxim Raginsky. Coordinate dual averaging for decentralized  
339 online optimization with nonseparable global objectives. *IEEE Transactions on Control of*  
340 *Network Systems*, 5(1):34–44, 2016.
- 341 [42] Aurélien Bellet, Rachid Guerraoui, Mahsa Taziki, and Marc Tommasi. Personalized and private  
342 peer-to-peer machine learning. *arXiv preprint arXiv:1705.08435*, 2017.
- 343 [43] Milad Nasr, Reza Shokri, and Amir Houmansadr. Comprehensive privacy analysis of deep  
344 learning: Stand-alone and federated learning under passive and active white-box inference  
345 attacks. *arXiv preprint arXiv:1812.00910*, 2018.
- 346 [44] Angelia Nedić and Alex Olshevsky. Distributed optimization over time-varying directed graphs.  
347 *IEEE Transactions on Automatic Control*, 60(3):601–615, 2014.
- 348 [45] Angelia Nedić and Alex Olshevsky. Stochastic gradient-push for strongly convex functions on  
349 time-varying directed graphs. *IEEE Transactions on Automatic Control*, 61(12):3936–3947,  
350 2016.
- 351 [46] Mahmoud Assran and Michael Rabbat. Asynchronous subgradient-push. *arXiv preprint*  
352 *arXiv:1803.08950*, 2018.
- 353 [47] Mahmoud Assran, Nicolas Loizou, Nicolas Ballas, and Michael Rabbat. Stochastic gradient  
354 push for distributed deep learning. *arXiv preprint arXiv:1811.10792*, 2018.

# Appendix

355

## 356 6 Related Work

357 The concept of *federated learning* was first proposed in [16], which advocates a novel learning setting  
358 that learns a shared model by aggregating locally-computed gradient updates without centralizing  
359 distributed data on devices. Early examples of research into federated learning also include [17, 18],  
360 and a widespread blog article posted by Google AI [19]. To address both statistical and system  
361 challenges, [20] and [21] propose a multi-task learning framework for federated learning and its  
362 related optimization algorithm, which extends early works SDCA [22, 23, 24] and COCOA [25,  
363 26, 27] to the federated learning setting. Among these optimization methods, *Federated Averaging*  
364 (FedAvg), proposed by [16], beats conventional synchronized mini-batch SGD with respect to  
365 communication rounds, and it converges on non-IID and unbalanced data. Recent rigorous theoretical  
366 analyses [28, 29, 30, 31] show that FedAvg is a special case of averaging periodic SGD (also  
367 called “local SGD”) which allows nodes to perform local updates and infrequent synchronization  
368 between them to communicate less while converging quickly. However, they cannot be applied to the  
369 single-sided trust network (asymmetric topology matrix).

370 Decentralized learning is a typical parallel strategy in which each worker is only required to com-  
371 municate with its neighbors, meaning that the communication bottleneck (in the parameter server)  
372 is removed. It has already been proven that decentralized learning can outperform the traditional  
373 centralized learning when the worker number is comparably large under a poor network condition  
374 [32]. There are two main types of decentralized learning algorithms: fixed network topology [33],  
375 and time-varying [34, 35] during training. [36, 37] show that the decentralized SGD can converge  
376 with a comparable convergence rate to the centralized algorithm with less communication to make  
377 large-scale model training feasible. [38] provides a systematic analysis of decentralized learning.

378 Online learning has been studied for decades. It is well known that the lower bounds of online  
379 optimization methods are  $\mathcal{O}(\sqrt{T})$  and  $\mathcal{O}(\log T)$  for convex and strongly convex loss functions,  
380 respectively [15, 8]. In recent years, due to the increasing volume of data, distributed online learning,  
381 especially decentralized methods, has attracted much attention. Examples of these works include  
382 [39, 40, 41]. Notably, [14] shares a similar problem definition and theoretical result as our paper.  
383 However, single-sided communication is not allowed in their setting, restricting their results.

## 384 7 Discussion on Privacy Protection

385 Although the main focus of this work is to provide a solution to deal with the practical scenario  
386 with asymmetric social networks, our proposed algorithm also has several advantages concerning  
387 privacy protection. First, as we have mentioned, OPS runs in a decentralized way and exchanges  
388 models instead of gradients or training samples, which is already proven effective for reducing the  
389 risk of privacy leakage [42]. Second, OPS runs in a decentralized and asymmetric fashion. These  
390 properties create difficulties for many attacking methods, such as [43]. In order to infer the data of  
391 other clients, the attacker needs to observe the reactions of other nodes after the attack is injected,  
392 which is impossible when the connections are single-sided. Even though the attack will spread among  
393 the whole network and finally return to the attacker, it is still hard for the attacker to distinguish  
394 whether the information he receives from its neighbors is already affected by the attack or not, since  
395 he is unaware of the global topology.

## 396 8 Extra Experiments

### 397 8.1 Evaluation on Different Network Sizes

398 Figure 2 summarizes the evaluation of OPS in different network sizes (in the *SUSY* dataset, 128,  
399 256, 512, 1024 are set, while in the *Room Occupancy* dataset, the network size is set to 10, 16, and  
400 20). The upper-bound neighbor number is aligned to the same value among different network sizes  
401 to isolate its impact. As we can see, in every dataset, the average loss (regret) curve in different

(a) stochastic=100%    (b) stochastic=50%    (c) stochastic=100%    (d) stochastic=50%

Figure 2: Evaluation on different network sizes

(a) stochastic=100%    (b) stochastic=50%    (c) stochastic=100%    (d) stochastic=50%

Figure 3: Evaluation on different network densities

402 network sizes is close on a small scale. These observations demonstrate OPS is robust to the network  
 403 size. Furthermore, the average loss (regret) is smaller in larger network size (i.e., the curve of the  
 404  $n = 1024$  network size is lower than others), which also demonstrates that more stochastic samples  
 405 provided by more nodes can naturally accelerate the convergence.

## 406 8.2 Evaluation on Network Density

407 We also evaluate the performance of OPS in different network densities. We set the network size  
 408 to 512 and 20 for SUSY and Room Occupancy dataset respectively. Network density is defined as  
 409 the ratio of the upper-bound random neighbor number per node to the size of the network (e.g., if  
 410 the ratio is 0.5 in SUSY it means 256 is set as the upper-bound neighbor number for each node).  
 411 We can see from Figure 3 that as the network density increased, the average loss (regret) decreased.  
 412 This observation also proves that our proposed OPS algorithm can work well in different network  
 413 densities, and can gain more benefits from a denser row stochastic matrix. This benefit can also be  
 414 understood intuitively: in a federated learning network, a user's model performance will improve if it  
 415 communicates with more users.

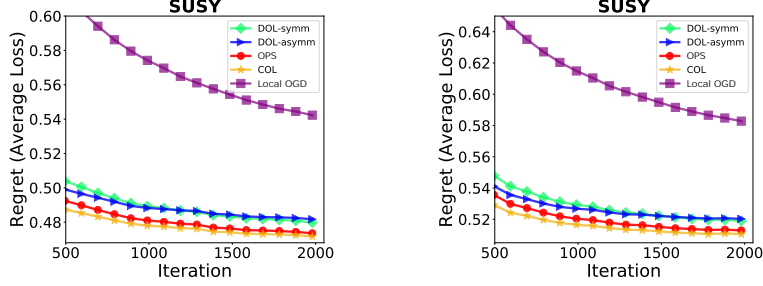
## 416 8.3 Comparison with local online gradient descent

417 To justify the necessity of communication, we also compare OPS with the local online gradient  
 418 descent (local OGD), where every node trains a local model without communicating with others. We  
 419 run experiments in different ratios of the adversary and stochastic components based on settings in  
 420 Figure 2. As we can see in Figure 4, we empirically prove that communication does have benefits in  
 421 reducing regret. Moreover, as the ratio of the stochastic components increased, the regret of OPS  
 422 decreases further. This also empirically proves that the stochastic component can benefit from the  
 423 communication while the adversarial component does not.

## 424 9 Proof to Theorem 2 and Theorem 3

425 Notations: Below we use the following notation in our proof

- 426 •  $r F_t(X_t) := \frac{1}{h} \sum_{i=1}^h r F_{1,t}(x_t^{(1)}); \dots; \frac{1}{h} \sum_{i=1}^h r F_{n,t}(x_t^{(n)})$
- 427 •  $X_t := x_t^{(1)}; x_t^{(2)}; \dots; x_t^{(n)}$
- 428 •  $G_t := \frac{1}{h} \sum_{i=1}^h f_{1,t}(x_t^1; \eta_t^1); \dots; \frac{1}{h} \sum_{i=1}^h f_{n,t}(x_t^n; \eta_t^n)$



(a) stochastic=100%  
( $n = 128, \#Neighbors=32$ )

(b) stochastic=50%  
( $n = 128, \#Neighbors=32$ )

Figure 4: Comparison between OPS and Local OGD.

429 Here we first present the proof Theorem 2, then we will present some key lemmas along with the  
430 proof of Theorem 3.

431 The following theorem is the key to prove Theorem 2:

432 **Theorem 4.** *For the online push-sum algorithm with step size  $\gamma > 0$ , it holds that*

$$\mathcal{R}_T \leq G^2 T n \gamma C_1 + \sigma^2 T \gamma (1 + n C_2) + \frac{n R^2}{2\gamma}, \quad (4)$$

433 where

$$C_1 := \frac{8Cq}{\delta_{\min}(1-q)} + 1, \quad C_2 := \frac{2Cq}{\delta_{\min}(1-q)},$$

434 and  $C$ ,  $q$  and  $\delta_{\min}$  are some constants defined in later lemmas.

435 *Proof.* Since the loss function  $f_{i,t}(\cdot)$  is assumed to be convex, which leads to

$$\begin{aligned} & \mathbb{E}_t \sum_{i=1}^n f_{i,t}(\mathbf{x}_t^{(i)}; \boldsymbol{\xi}_t^{(i)}) - n F_t(\mathbf{x}^*) \\ &= \mathbb{E}_t \sum_{i=1}^n \left( f_{i,t}(\mathbf{x}_t^{(i)}; \boldsymbol{\xi}_t^{(i)}) - f_{i,t}(\mathbf{x}^*; \boldsymbol{\xi}_t^{(i)}) \right) \\ &\leq \mathbb{E}_t \sum_{i=1}^n \left\langle \nabla f_{i,t}(\mathbf{x}_t^{(i)}; \boldsymbol{\xi}_t^{(i)}), \mathbf{x}_t^{(i)} - \mathbf{x}^* \right\rangle \\ &= \underbrace{\mathbb{E}_t \sum_{i=1}^n \left\langle \nabla f_{i,t}(\mathbf{x}_t^{(i)}; \boldsymbol{\xi}_t^{(i)}), \mathbf{x}_t^{(i)} - \bar{\mathbf{z}}_t \right\rangle}_{:= I_{1t}} + \underbrace{\mathbb{E}_t \sum_{i=1}^n \left\langle \nabla f_{i,t}(\mathbf{x}_t^{(i)}; \boldsymbol{\xi}_t^{(i)}), \bar{\mathbf{z}}_t - \mathbf{x}^* \right\rangle}_{:= I_{2t}}. \end{aligned}$$

436 For  $I_{2t}$ , we have

$$\begin{aligned}
& \mathbb{E}_t \sum_{i=1}^n \left\langle \nabla f_{i,t} \left( \mathbf{x}_t^{(i)}; \boldsymbol{\xi}_t^{(i)} \right), \bar{\mathbf{z}}_t - \mathbf{x}^* \right\rangle \\
&= \frac{n}{\gamma} \mathbb{E}_t \left\langle \frac{\gamma}{n} \sum_{i=1}^n \nabla f_{i,t} \left( \mathbf{x}_t^{(i)}; \boldsymbol{\xi}_t^{(i)} \right), \bar{\mathbf{z}}_t - \mathbf{x}^* \right\rangle \\
&= \frac{n}{2\gamma} \mathbb{E}_t \left( \left\| \frac{\gamma}{n} \sum_{i=1}^n \nabla f_{i,t} \left( \mathbf{x}_t^{(i)}; \boldsymbol{\xi}_t^{(i)} \right) \right\|^2 + \|\bar{\mathbf{z}}_t - \mathbf{x}^*\|^2 - \left\| \bar{\mathbf{z}}_t - \mathbf{x}^* - \frac{\gamma}{n} \sum_{i=1}^n \nabla f_{i,t} \left( \mathbf{x}_t^{(i)}; \boldsymbol{\xi}_t^{(i)} \right) \right\|^2 \right) \\
&= \frac{n}{2\gamma} \mathbb{E}_t \left( \left\| \frac{\gamma}{n} \sum_{i=1}^n \nabla f_{i,t} \left( \mathbf{x}_t^{(i)}; \boldsymbol{\xi}_t^{(i)} \right) \right\|^2 + \|\bar{\mathbf{z}}_t - \mathbf{x}^*\|^2 - \|\bar{\mathbf{z}}_{t+1} - \mathbf{x}^*\|^2 \right) \\
&\leq \frac{n}{2\gamma} \mathbb{E}_t \left( \gamma^2 G^2 + \frac{\gamma^2 \sigma^2}{n} + \|\bar{\mathbf{z}}_t - \mathbf{x}^*\|^2 - \|\bar{\mathbf{z}}_{t+1} - \mathbf{x}^*\|^2 \right)
\end{aligned}$$

437 Notice that for COL, we have  $I_{1t} = 0$  because  $\mathbf{x}_t^{(i)} = \bar{\mathbf{z}}_t$ . So for DOL, in order to bound  $I_{1t}$ , we need

438 to bound the difference  $\left\| \mathbf{x}_t^{(i)} - \bar{\mathbf{z}}_t \right\|$  (using Lemma 8).

$$\begin{aligned}
& \mathbb{E}_t \sum_{i=1}^n \left\langle \nabla f_{i,t} \left( \mathbf{x}_t^{(i)}; \boldsymbol{\xi}_t^{(i)} \right), \mathbf{x}_t^{(i)} - \bar{\mathbf{z}}_t \right\rangle \\
&= \mathbb{E}_t \sum_{i=1}^n \left\langle \nabla F_{i,t}(\mathbf{x}_t^{(i)}), \mathbf{x}_t^{(i)} - \bar{\mathbf{z}}_t \right\rangle \\
&\leq \mathbb{E}_t \sum_{i=1}^n \left( \alpha \left\| \nabla F_{i,t} \left( \mathbf{x}_t^{(i)} \right) \right\|^2 + \frac{1}{\alpha} \left\| \mathbf{x}_t^{(i)} - \bar{\mathbf{z}}_t \right\|^2 \right).
\end{aligned}$$

439 Summing up the inequality above from  $t = 1$  to  $t = T$ , we get

$$\begin{aligned}
& \sum_{t=1}^T \mathbb{E}_t \sum_{i=1}^n \left\langle \nabla f_{i,t} \left( \mathbf{x}_t^{(i)}; \boldsymbol{\xi}_t^{(i)} \right), \mathbf{x}_t^{(i)} - \bar{\mathbf{z}}_t \right\rangle \\
&= \sum_{t=1}^T \mathbb{E}_t \sum_{i=1}^n \left\langle \nabla F_{i,t} \left( \mathbf{x}_t^{(i)} \right), \mathbf{x}_t^{(i)} - \bar{\mathbf{z}}_t \right\rangle \\
&\leq \sum_{t=1}^T \mathbb{E}_t \sum_{i=1}^n \left( \alpha \left\| \nabla F_{i,t} \left( \mathbf{x}_t^{(i)} \right) \right\|^2 + \frac{1}{\alpha} \left\| \mathbf{x}_t^{(i)} - \bar{\mathbf{z}}_t \right\|^2 \right) \\
&= \sum_{t=1}^T \left( \alpha \mathbb{E}_t \left\| \nabla F_t(\mathbf{X}_t) \right\|_F^2 + \frac{1}{\alpha} \mathbb{E}_t \left\| \mathbf{X}_t - \bar{\mathbf{z}}_t \right\|_F^2 \right) \\
&\leq \alpha \sum_{t=1}^T \mathbb{E}_t \left\| \nabla F_t(\mathbf{X}_t) \right\|_F^2 + \frac{4\gamma^2 C^2 q^2}{\alpha \delta_{\min}^2 (1-q)^2} \sum_{t=1}^T \mathbb{E}_t \left\| \mathbf{G}_t \right\|_F^2 \\
&\leq \alpha \sum_{t=1}^T \mathbb{E}_t \left\| \nabla F_t(\mathbf{X}_t) \right\|_F^2 + \frac{4\gamma^2 C^2 q^2}{\alpha \delta_{\min}^2 (1-q)^2} \sum_{t=1}^T \left( \mathbb{E}_t \left\| \nabla F_t(\mathbf{X}_t) \right\|_F^2 + n\sigma^2 \right).
\end{aligned}$$

440 Choosing  $\alpha = \frac{2\gamma C q}{\delta_{\min}(1-q)}$ , we have

$$\sum_{t=1}^T \mathbb{E}_t \sum_{i=1}^n \left\langle \nabla f_{i,t} \left( \mathbf{x}_t^{(i)}; \boldsymbol{\xi}_t^{(i)} \right), \mathbf{x}_t^{(i)} - \bar{\mathbf{z}}_t \right\rangle \leq \frac{8n\gamma C T q G^2}{\delta_{\min}(1-q)} + \frac{2n\gamma C q \sigma^2 T}{\delta_{\min}(1-q)}$$

441 So we have

$$\begin{aligned}
& \sum_{t=1}^T \mathbb{E}_t \sum_{i=1}^n f_{i,t} \left( \mathbf{z}_t^{(i)}; \boldsymbol{\xi}_t^{(i)} \right) - nF(\mathbf{x}^*) \\
& \leq \frac{8n\gamma CTqG^2}{\delta_{\min}(1-q)} + \frac{2\gamma Cq\sigma^2 T}{\delta_{\min}(1-q)} + \frac{n}{2n\gamma} \sum_{t=1}^T \left( \gamma^2 G^2 + \frac{\gamma^2 \sigma^2}{n} + \mathbb{E}_t \|\bar{\mathbf{z}}_t - \mathbf{x}^*\|^2 - \mathbb{E}_t \|\bar{\mathbf{z}}_{t+1} - \mathbf{x}^*\|^2 \right) \\
& \leq G^2 T n \gamma \left( \frac{8Cq}{\delta_{\min}(1-q)} + 1 \right) + \sigma^2 T \gamma \left( 1 + \frac{2nCq}{\delta_{\min}(1-q)} \right) + \frac{n}{2\gamma} \sum_{t=1}^T \left( \mathbb{E}_t \|\bar{\mathbf{z}}_t - \mathbf{x}^*\|^2 - \mathbb{E}_t \|\bar{\mathbf{z}}_{t+1} - \mathbf{x}^*\|^2 \right) \\
& \leq G^2 T n \gamma \left( \frac{8Cq}{\delta_{\min}(1-q)} + 1 \right) + \sigma^2 T \gamma \left( 1 + \frac{2nCq}{\delta_{\min}(1-q)} \right) + \frac{nR^2}{2\gamma} \\
& = C_1 n G^2 T \gamma + (1 + nC_2) \sigma^2 T \gamma + \frac{nR^2}{2\gamma}.
\end{aligned}$$

442 Notice that Theorem 2 can be easily verified by setting  $\gamma = \frac{\sqrt{n}R}{\sqrt{(1+nC_2)\sigma^2 + \sqrt{n}C_1 G^2 T}}$ . □

443 Next, we will present two lemmas for our proof of Lemma 8. The proofs of following two lemmas  
444 can be found in existing literature [44, 45, 46, 47].

445 **Lemma 5.** *Under the Assumption 1, there exists a constant  $\delta_{\min} > 0$  such that for any  $t$ , the following*  
446 *holds*

$$\sum_{j=1}^n [\mathbf{W}^{t\top} \mathbf{W}^{t\top} \dots \mathbf{W}^{0\top}]_{ij} \geq \delta_{\min} \geq \frac{1}{n^n}, \forall i \quad (5)$$

447 where  $\mathbf{W}^t$  is a row stochastic matrix.

448 **Lemma 6.** *Under the Assumption 1, for any  $t$ , there always exists a stochastic vector  $\psi(t)$  and two*  
449 *constants  $C = 4$  and  $q = 1 - n^{-n} < 1$  such that for any  $s$  satisfying  $s \leq t$ , the following inequality*  
450 *holds*

$$|[\mathbf{W}^{t\top} \mathbf{W}^{t\top} \dots \mathbf{W}^{s+1\top} \mathbf{W}^{s\top}]_{ij} - \psi_i(t)| \leq Cq^{t-s}, \forall i, j$$

451 where  $\mathbf{W}^t$  is a row stochastic matrix, and  $\psi(t)$  is a vector with  $\psi_i(t)$  being its  $i$ -th entry.

452 **Lemma 7.** *Given two non-negative sequences  $\{a_t\}_{t=1}^\infty$  and  $\{b_t\}_{t=1}^\infty$  that satisfying*

$$a_t = \sum_{s=1}^t \rho^{t-s} b_s, \quad (6)$$

453 with  $\rho \in [0, 1)$ , we have

$$D_k := \sum_{t=1}^k a_t^2 \leq \frac{1}{(1-\rho)^2} \sum_{s=1}^k b_s^2.$$

454 *Proof.* From the definition, we have

$$\begin{aligned}
S_k &= \sum_{t=1}^k \sum_{s=1}^t \rho^{t-s} b_s = \sum_{s=1}^k \sum_{t=s}^k \rho^{t-s} b_s = \sum_{s=1}^k \sum_{t=0}^{k-s} \rho^t b_s \leq \sum_{s=1}^k \frac{b_s}{1-\rho}, \quad (7) \\
D_k &= \sum_{t=1}^k \sum_{s=1}^t \rho^{t-s} b_s \sum_{r=1}^t \rho^{t-r} b_r \\
&= \sum_{t=1}^k \sum_{s=1}^t \sum_{r=1}^t \rho^{2t-s-r} b_s b_r \\
&\leq \sum_{t=1}^k \sum_{s=1}^t \sum_{r=1}^t \rho^{2t-s-r} \frac{b_s^2 + b_r^2}{2} \\
&= \sum_{t=1}^k \sum_{s=1}^t \sum_{r=1}^t \rho^{2t-s-r} b_s^2 \\
&\leq \frac{1}{1-\rho} \sum_{t=1}^k \sum_{s=1}^t \rho^{t-s} b_s^2 \\
&\leq \frac{1}{(1-\rho)^2} \sum_{s=1}^k b_s^2.
\end{aligned}$$

455

□

456 Based on the above three lemmas, we can obtain the following lemma.

457 **Lemma 8.** *Under the Assumption 1, the updating rule of Algorithm 1 leads to the following inequality*

$$\sum_i^n \sum_{t=0}^T \left\| \mathbf{x}_{t+1}^{(i)} - \bar{\mathbf{z}}_{t+1} \right\|_2^2 \leq \frac{4\gamma^2 C^2 q^2}{\delta_{\min}^2 (1-q)^2} \sum_{s=0}^t \|\mathbf{G}_s\|_F^2,$$

458 where  $\gamma$  is the step size, and  $C = 4$ ,  $\delta_{\min} \geq n^{-n}$ ,  $q = 1 - n^{-n}$  are constants.  $\mathbf{G}_s$  is the matrix for  
459 the stochastic gradient at time  $s$  (e.g., the  $i$ -th column is the stochastic gradient vector on node  $i$  at  
460 time  $s$ ).

461 *Proof.* The updating rule of OPS can be formulated as

$$\begin{aligned}
\mathbf{Z}_{t+1} &= (\mathbf{Z}_t - \gamma \mathbf{G}_t) \mathbf{W} \\
\omega_{t+1} &= \mathbf{W}^\top \omega_t \\
\mathbf{X}_{t+1} &= \mathbf{Z}_{t+1} [\text{diag}(\omega_{t+1})]^{-1}
\end{aligned}$$

462 where  $\mathbf{W}$  is a row stochastic matrix.  $\mathbf{X}_t = [\mathbf{x}_t^{(1)}, \mathbf{x}_t^{(2)}, \dots, \mathbf{x}_t^{(n)}]$  is a matrix whose each column is  
463  $\mathbf{x}_t^{(i)}$ .  $\mathbf{G}_t$  is the matrix of gradient, whose each column is the stochastic gradient at  $\mathbf{z}_t^{(i)}$  on node  $i$ .  
464  $\mathbf{Z}_t = [\mathbf{z}_t^{(1)}, \dots, \mathbf{z}_t^{(n)}]$  is the matrix whose each column is  $\mathbf{z}_t^{(i)}$ .

465 Assuming  $\mathbf{X}_0 = \mathbf{O}$  and  $\omega_0 = 1$ , then we have

$$\mathbf{Z}_{t+1} = (\mathbf{Z}_t - \gamma \mathbf{G}_t) \mathbf{W} = \dots = -\gamma \sum_{s=0}^t \mathbf{G}_s \mathbf{W}^{t-s+1}, \quad (8)$$

$$\bar{\mathbf{z}}_{t+1} = \bar{\mathbf{z}}_t - \gamma \bar{\mathbf{g}}_t = \dots = -\sum_{s=0}^t \gamma \bar{\mathbf{g}}_s, \quad (9)$$

$$\omega_{t+1} = \mathbf{W}^{t+1\top} \omega_0, \quad (10)$$

466 where  $\bar{\mathbf{x}}_t = \mathbf{X}_t \mathbf{1}$  is the average of all variables on the  $n$  nodes, and  $\bar{\mathbf{g}}_t = \mathbf{G}_t \mathbf{1}$  is the averaged gradient.  
467 We have  $\mathbf{W} \mathbf{1} = \mathbf{1}$  since  $\mathbf{W}$  is a row stochastic matrix.

468 For  $\omega_{t+1}$ , according to Lemma 6, we decompose it as follows

$$\omega_{t+1} = \mathbf{W}^{t+1\top} \omega_0 = [\mathbf{W}^{t+1\top} - \psi(t)\mathbf{1}^\top] \omega_0 + \psi(t)\mathbf{1}^\top \omega_0 = [\mathbf{W}^{t+1\top} - \psi(t)\mathbf{1}^\top] \mathbf{1} + n\psi(t), \quad (11)$$

469 since  $\omega_0 = \mathbf{1}$ .

470 On the other hand, according to Lemma 5, we also have

$$\omega_{t+1}^{(i)} = [\mathbf{W}^{t+1\top} \mathbf{1}]^\top \mathbf{e}_i = \sum_{j=1}^n [\mathbf{W}^{t+1\top}]_{ij} \geq n\delta_{\min}, \quad (12)$$

471 where  $\mathbf{e}_i$  is a vector with only the  $i$ -th entry being 1 and 0 for others.

472 We need to further bound the following term

$$\begin{aligned} \left\| \mathbf{x}_{t+1}^{(i)} - \bar{\mathbf{z}}_{t+1} \right\| &= \gamma \left\| \frac{\mathbf{z}_{t+1}^{(i)}}{\omega_{t+1}^{(i)}} - \bar{\mathbf{z}}_{t+1} \right\| \\ &= \gamma \left\| \sum_{s=0}^t \left( \frac{\mathbf{G}_s \mathbf{W}^{t-s+1} \mathbf{e}_i}{\mathbf{1}^\top \mathbf{W}^{t+1} \mathbf{e}_i} - \frac{\mathbf{G}_s \mathbf{1}}{n} \right) \right\| \\ &= \gamma \left\| \sum_{s=0}^t \frac{n\mathbf{G}_s \mathbf{W}^{t-s+1} \mathbf{e}_i - \mathbf{G}_s \mathbf{1} \mathbf{1}^\top \mathbf{W}^{t+1} \mathbf{e}_i}{n\omega_{t+1}^{(i)}} \right\|, \end{aligned}$$

473 where the second equality is by (8), (9), and (10). We turn to bound the following term

$$\begin{aligned} &\left\| \sum_{s=0}^t \frac{n\mathbf{G}_s \mathbf{W}^{t-s+1} \mathbf{e}_i - \mathbf{G}_s \mathbf{1} \mathbf{1}^\top \mathbf{W}^{t+1} \mathbf{e}_i}{n\omega_{t+1}^{(i)}} \right\| \\ &\leq \frac{1}{n^2 \delta_{\min}^2} \left\| \sum_{s=0}^t (n\mathbf{G}_s \mathbf{W}^{t-s+1} \mathbf{e}_i - \mathbf{G}_s \mathbf{1} \mathbf{1}^\top \mathbf{W}^{t+1} \mathbf{e}_i) \right\|, \end{aligned}$$

474 where the first inequality is according to (12). Therefore, combining the results above, we can have

$$\begin{aligned} \sum_{i=1}^n \left\| \mathbf{x}_{t+1}^{(i)} - \bar{\mathbf{z}}_{t+1} \right\|_2^2 &\leq \frac{\gamma^2}{n^4 \delta_{\min}^2} \sum_{i=1}^n \left\| \sum_{s=0}^t (n\mathbf{G}_s \mathbf{W}^{t-s+1} \mathbf{e}_i - \mathbf{G}_s \mathbf{1} \mathbf{1}^\top \mathbf{W}^{t+1} \mathbf{e}_i) \right\|_2^2 \\ &\leq \frac{\gamma^2}{n^4 \delta_{\min}^2} \left\| \sum_{s=0}^t (n\mathbf{G}_s \mathbf{W}^{t-s+1} - \mathbf{G}_s \mathbf{1} \mathbf{1}^\top \mathbf{W}^{t+1}) \right\|_F^2 \end{aligned}$$

475 where the second inequality is due to  $\sum_{i=1}^n \|\mathbf{A} \mathbf{e}_i\|_2^2 = \|\mathbf{A}\|_F^2$ .



476 It remains to bound the following term

$$\begin{aligned}
& \left\| \sum_{s=0}^t (n\mathbf{G}_s \mathbf{W}^{t-s+1} - \mathbf{G}_s \mathbf{1} \mathbf{1}^\top \mathbf{W}^{t+1}) \right\|_F^2 \\
&= \left\| \sum_{s=0}^t (n\mathbf{G}_s \mathbf{W}^{t-s+1} - \mathbf{G}_s \mathbf{1} [1^\top (\mathbf{W}^{t+1} - \psi(t) \mathbf{1}^\top)^\top + n\psi(t)^\top]) \right\|_F^2 \\
&= \left\| \sum_{s=0}^t (n\mathbf{G}_s [\mathbf{W}^{t-s+1} - \mathbf{1} \psi(t)^\top] - \mathbf{G}_s \mathbf{1} \mathbf{1}^\top [\mathbf{W}^{t+1} - \mathbf{1} \psi(t)^\top]) \right\|_F^2 \\
&\leq \left( \sum_{s=0}^t \|n\mathbf{G}_s [\mathbf{W}^{t-s+1} - \mathbf{1} \psi(t)^\top]\|_F + \sum_{s=0}^t \|\mathbf{G}_s \mathbf{1} \mathbf{1}^\top [\mathbf{W}^{t+1} - \mathbf{1} \psi(t)^\top]\|_F \right)^2 \\
&\leq \left( n \sum_{s=0}^t \|\mathbf{G}_s\|_F \|[\mathbf{W}^{t-s+1} - \mathbf{1} \psi(t)^\top]\|_F + \sum_{s=0}^t \|\mathbf{G}_s\|_F \|\mathbf{1} \mathbf{1}^\top\|_F \|[\mathbf{W}^{t+1} - \mathbf{1} \psi(t)^\top]\|_F \right)^2 \\
&\leq n^2 \left( \sum_{s=0}^t \|\mathbf{G}_s\|_F \|[\mathbf{W}^{t-s+1} - \mathbf{1} \psi(t)^\top]\|_F + \sum_{s=0}^t \|\mathbf{G}_s\|_F \|[\mathbf{W}^{t+1} - \mathbf{1} \psi(t)^\top]\|_F \right)^2 \\
&\leq n^2 \left( \sum_{s=0}^t nCq^{t-s+1} \|\mathbf{G}_s\|_F + \sum_{s=0}^t nCq^{t+1} \|\mathbf{G}_s\|_F \right)^2 \\
&\leq 4n^4 C^2 q^2 \left( \sum_{s=0}^t q^{t-s} \|\mathbf{G}_s\|_F \right)^2
\end{aligned}$$

477 where the third inequality is due to  $\|\mathbf{1} \mathbf{1}^\top\|_F = n$  and the fourth inequality is by Lemma 6 and the  
478 fact that  $\|\mathbf{A}\|_F \leq n \cdot \max_{i,j} |A_{ij}|$  if  $\mathbf{A} \in \mathbb{R}^{n \times n}$ .

479 Therefore, if we combining all the above inequalities together, we can obtain

$$\sum_{i=1}^n \left\| \mathbf{x}_{t+1}^{(i)} - \bar{\mathbf{z}}_{t+1} \right\|_2^2 \leq \frac{4\gamma^2 C^2 q^2}{\delta_{\min}^2} \left( \sum_{s=0}^t q^{t-s} \|\mathbf{G}_s\|_F \right)^2.$$

480 Using Lemma 7, we have

$$\sum_{t=0}^T \left( \sum_{s=0}^t q^{t-s} \|\mathbf{G}_s\|_F \right)^2 \leq \frac{1}{(1-q)^2} \sum_{t=0}^T \|\mathbf{G}_t\|_F^2,$$

481 which leads to

$$\sum_{t=0}^T \sum_{i=1}^n \left\| \mathbf{x}_{t+1}^{(i)} - \bar{\mathbf{z}}_{t+1} \right\|_2^2 \leq \frac{4\gamma^2 C^2 q^2}{\delta_{\min}^2 (1-q)^2} \sum_{t=0}^T \|\mathbf{G}_t\|_F^2,$$

482 which completes the proof.  $\square$

483 Actually, Theorem 3 is a corollary of Lemma 8 by setting  $\gamma$  as the appropriate value.